

FoMo-0D: A Foundation Model for Zero-shot Outlier Detection

Anonymous authors
Paper under double-blind review

Abstract

Outlier detection (OD) has a vast literature as it finds numerous real-world applications. Being an unsupervised task, model selection is a key bottleneck for OD without label supervision. Despite a long list of available OD algorithms with tunable hyperparameters, the lack of systematic approaches for unsupervised algorithm and hyperparameter selection limits their effective use in practice. In this paper, we present FoMo-0D, a pre-trained Foundation Model for zero/0-shot OD on tabular data, which bypasses the hurdle of model selection altogether. Capitalizing on synthetic pre-training, FoMo-0D can directly predict the (outlier/inlier) label of test samples without parameter fine-tuning. Even more importantly, it *requires no labeled data, and no additional training or hyperparameter tuning when given a new task*. Extensive experiments on 57 real-world datasets against 26 baselines show that FoMo-0D is highly competitive; outperforming the majority of the baselines with no statistically significant difference from the 2nd best method. Further, FoMo-0D is efficient in inference time requiring only 7.7 ms per sample on average, with at least 7x speed-up compared to previous methods. To facilitate future research, our implementations for data synthesis and pre-training as well as model checkpoints are openly available at <https://anonymous.4open.science/r/PFN40D>.

1 Introduction

Outlier detection (OD) in tabular data finds numerous applications in critical domains such as security, environmental monitoring, finance, and medicine, to name a few. This popularity brings along a large literature with plethora of detection algorithms to choose from given a new OD task. These algorithms, however, exhibit several hyperparameters (HPs) that need careful tuning (Ma et al., 2023). Since most OD tasks are unsupervised¹, what makes effective detection notoriously difficult is unsupervised model selection (both algorithm and HP selection) in the absence of labels.

While deep learning has revolutionized many areas of machine learning (ML), it is not quite the case for OD. This is mainly because compared to classical methods, deep OD models (Pang et al., 2021) have many more HPs to which the detection performance is sensitive (Ding et al., 2022), making model selection even more challenging. While the recent success of large foundation models, pre-trained on massive amounts of data, offers new opportunities for zero-shot OD, thus far the most notable progress has been in NLP and computer vision Brown et al. (2020); Touvron et al. (2023); Radford et al. (2021). This is thanks to the admirable quantity and quality of public text and image datasets. In comparison, public tabular OD benchmarks remain minuscule (Han et al., 2022; Zhao et al., 2021; Steinbuss & Böhm, 2021).

Recently, Prior-data Fitted Networks (PFNs) has marked a milestone in ML as a new approach to learning on tabular data (Müller et al., 2022). The core idea is to compute a posterior predictive distribution (PPD) for a test point given training data as context. To approximate the PPD, a Transformer (Vaswani et al., 2017) is pre-trained on a large set of synthetic datasets drawn from pre-defined data priors. At inference, the pre-trained PFN is fed with test samples along with some training samples as context for zero-shot prediction, requiring no parameter fine-tuning or model selection on new datasets. Variants of PFN are shown to match

¹While supervised OD exists, unsupervised setting is preferred in most domains to detect novel, emergent anomalies.

Table 1: p -values of the one-sided Wilcoxon signed rank test, comparing FoMo-0D (with $D = 100$) to **top 10 baselines** with default hyperparameters (HPs), and **top 4^{avg}** baselines⁶ with **avg.** performance over varying HPs (denoted w/ ^{avg}) over All (57) datasets, those (42) w/ $d \leq 100$ and (46) w/ $d \leq 500$ dimensions, and those (47) excluding NLP, CV datasets. FoMo-0D shows **no statistically significant difference from the 2nd best model** ($k\text{NN}$, w/ $p = 0.106$) over All datasets, and **none of the differences are significant when embedded, i.e., NLP and CV datasets are excluded**, while it is comparable to ($p > \alpha$) or significantly better than ($p > 1 - \alpha$) all 26 original + 4^{avg} baselines over datasets w/ $d \leq 100$ (aligned w/ pretraining where $D = 100$) as well as on datasets w/ $d \leq 500$ (generalizing beyond pretraining). We use underline and **bold** to indicate $p < \alpha$ and $p > 1 - \alpha$. Rank, avg.'ed over all 57 datasets by AUROC. (setting: $D = 100$, $R = 500$, train/inference context size=5K, w/ quantile transform, $\alpha = 0.05$) (See Tables 17.1&17.2 for full results.)

	FoMo-0D	DTE-NP	<u>k</u> NN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP ^{avg}	<u>k</u> NN ^{avg}	ICL ^{avg}	DTE-C ^{avg}
$d \leq 100$	-	0.415	0.700	0.949	0.953	0.970	0.971	0.996	0.876	0.980	0.978	0.752	0.860	0.958	1.000
$d \leq 500$	-	0.220	0.569	0.827	0.894	0.960	0.968	0.994	0.910	0.960	0.979	0.607	0.756	0.846	1.000
All	-	0.016	0.106	0.462	0.454	0.585	0.750	0.823	0.759	0.901	0.895	0.112	0.315	0.670	1.000
All\{NLP, CV\}	-	0.164	0.476	0.757	0.832	0.934	0.945	0.988	0.867	0.938	0.965	0.515	0.683	0.777	1.000
Rank(avg)	11.886	7.553	9.018	10.851	11.36	12.316	13.342	13.386	12.982	14.061	13.851	9.079	11.105	12.991	22.263

37 tree-based models in performance on small classification datasets (Hollmann et al., 2023) as well as time
 38 series forecasting (Dooley et al., 2023).

39 In this paper, we capitalize on these ideas and introduce FoMo-0D: a prior-data fitted Foundation Model for
 40 zero/0-shot OD. Once pre-trained on synthetic datasets, FoMo-0D unlocks zero-shot OD on a new dataset
 41 where the (unlabeled) input data is fed only as context. As such, FoMo-0D bypasses not only model (parameter)
 42 training, but more importantly, the nontrivial task of unsupervised model (algorithm and HP) selection
 43 without labeled data. Figure 1 illustrates the new FoMo-0D paradigm versus the typical OD setting. To our
 44 knowledge, FoMo-0D is the first pre-trained foundation model for tabular OD.

45 In designing FoMo-0D, we use Gaussian mixture models as a simple yet effective tabular data prior for inlier
 46 data distributions (Hollmann et al., 2023; Zhao et al., 2021), which are also employed to simulate outlier
 47 types common in the real world; namely, local and global subspace outliers (Steinbuss & Böhm, 2021). While
 48 the data prior can be extended to comprise more complex data distributions (Hollmann et al., 2023) (e.g.
 49 Bayesian Neural Networks (Neal, 2012) and Structural Causal Models (Pearl, 2009)), and additional outlier
 50 types can be included (e.g. dependency, contextual, etc.), as we show in the experiments, even with its
 51 relatively straightforward prior, FoMo-0D achieves remarkable performance: As shown in Table 1 FoMo-0D,
 52 which is pre-trained on datasets with $d \leq 100$ dimensions, shows no statistically significant difference from
 53 all 26 state-of-the-art baselines (all p -values > 0.4) on 42 benchmark datasets with dimensionality $d \leq 100$
 54 (aligned with pre-training), while our method consistently ranks among the top and outperforms a majority
 55 of the baselines with p -value > 0.95 . (See Appendix Tables 17.1&17.2 for full results.) The results remain
 56 consistent on (46) benchmarks with $d \leq 500$ dimensions. FoMo-0D is also competitive across all (57) datasets,
 57 effectively generalizing beyond its pre-training distributions, with no statistically significant difference from
 58 the 2nd best baseline. When intrinsically non-tabular datasets are removed (i.e., datasets from NLP, CV
 59 domains embedded with pre-trained encoders), there is no statistical evidence to suggest a performance
 60 difference between FoMo-0D and any baseline on the remaining (47) datasets. Further, FoMo-0D takes a mere
 61 7.7 ms to infer a test sample on average with no extra training or tuning overhead on the new dataset. We
 62 summarize the main contributions of our work as follows.

63 • **A Foundation Model for Tabular OD:** We present FoMo-0D, *the first foundation model for zero-*
 64 *shot OD* on unseen tabular datasets, with no additional training or hyperparameter tuning, backed by
 65 Transformer-based in-context learning (ICL), synthetic data pre-training, and feed-forward inference.

66 • **Model Selection Made Obsolete:** FoMo-0D is designed for zero-shot inference given a new dataset,
 67 fully abolishing not only model training on the new dataset, but also the notorious task of algorithm
 68 selection and hyperparameter tuning in the absence of labeled data.

- **Scalable Pre-training:** To enable pre-training on many large datasets, we propose (i) a new mechanism to reduce sample-to-sample attention from quadratic to linear time—enabling larger datasets, as well as (ii) on-the-fly data synthesis through data transformations—enabling more diverse datasets in less time.
- **Fast OD at Inference:** Given a new dataset, FoMo-0D bypasses both model training and selection, both of which can be slow for modern deep OD models with many hyperparameters. Rather, it takes fraction of a second to label a test point through a single forward pass. Such speedy inference also unlocks the potential for deploying FoMo-0D in real time on data streams.
- **Effectiveness:** On a large benchmark of **57** datasets (Han et al., 2022) from diverse domains and against **26** baselines ranging from classical to modern OD models (Livernoche et al., 2024), FoMo-0D outperforms the majority of the baselines, with no statistical evidence for performance difference from the *2nd* best baseline, while operating fully zero-shot on real-world datasets out-of-the-box.
- **New directions:** As the first foundation model for OD, FoMo-0D presents a paradigm shift in how to perform OD in practice, while offering new directions to explore as well as open questions to investigate. From the algorithmic perspective, how can we understand what (OD) algorithm, if any, has the Transformer- and ICL-based FoMo-0D learned? How can we interpret (mechanistically or otherwise) how the pretraining achieves zero-shot and out-of-distribution (OOD) generalization? From a data perspective, what prior distributions for pretraining are suitable for downstream real-world tasks? What is the role of prior data diversity and complexity in the generalization ability of the model? In summary, FoMo-0D paves the path towards a better understanding of ICL as well as building more powerful future foundation models for OD.

88 2 Problem and Preliminaries

89 2.1 Outlier Detection Problem and Setting

90 Outlier detection (OD) methods can be categorized based on the availability of labeled data. In supervised
 91 OD, the task is similar to binary classification with imbalanced classes (as outliers typically make up only a
 92 small portion of the overall data). The more difficult unsupervised setting assumes that the “contaminated”
 93 training data contains both inliers and outliers, but without any labels. This is also the transductive setting
 94 where training and test data are the same. The one-class setting lies between these two extremes, where the
 95 “clean” training data consists only of inliers, while the unlabeled test data contains the outliers to be detected.
 96 Here, training and test sets are disjoint and thus the setting is inductive. Note that there exist **no labeled**
 97 **outliers** in both settings, making model selection challenging.

98 We remark that some OD literature refers to the latter setting as semi-supervised OD, which is a misnomer
 99 from the supervised ML perspective where semi-supervised classification assumes the presence of some labeled
 100 instances from **all** classes in the training data. In the rest of text, we adopt the terminology unsupervised OD
 101 for both settings, and specify “clean” inlier-only versus “contaminated” mixed training data to differentiate
 102 them. Then, our work considers the unsupervised OD problem under “clean” inlier-only training data.

103 Formally, let $\mathcal{D}_{\text{in}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ denote the inlier-only input data $\mathbf{x}_i \in \mathbb{R}^d$, where $y_i = 0 \forall i \in [n]$,
 104 and $\mathcal{D}_{\text{test}}$ depicts the unlabeled test data comprising both inliers and outliers. Note that $\mathcal{D}_{\text{in}} \cap \mathcal{D}_{\text{test}} = \emptyset$, i.e.
 105 **train/test split is disjoint**. The task is to assign labels to $\mathbf{x}_i \in \mathcal{D}_{\text{test}}$ given the inlier-only input \mathcal{D}_{in} .

106 2.2 Background on Prior-data Fitted Networks

107 **Posterior Predictive Distribution (PPD):** In the Bayesian framework for supervised learning, the prior
 108 defines a hypothesis space Φ which expresses our beliefs about the data distribution before seeing any data.
 109 Each hypothesis $\phi \in \Phi$ describes a mechanism by which the data is generated. The posterior predictive
 110 distribution $p(\cdot | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ provides a framework for making prediction on new, unseen test data \mathbf{x}_{test} ,
 111 conditioned on observed training data $\mathcal{D}_{\text{train}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. Based on Bayes’ Theorem, the PPD
 112 can be derived by the integration over the space of hypotheses Φ :

$$p(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}}) = \int_{\Phi} p(y_{\text{test}} | \mathbf{x}_{\text{test}}, \phi) p(\mathcal{D}_{\text{train}} | \phi) p(\phi) d\phi, \quad (1)$$

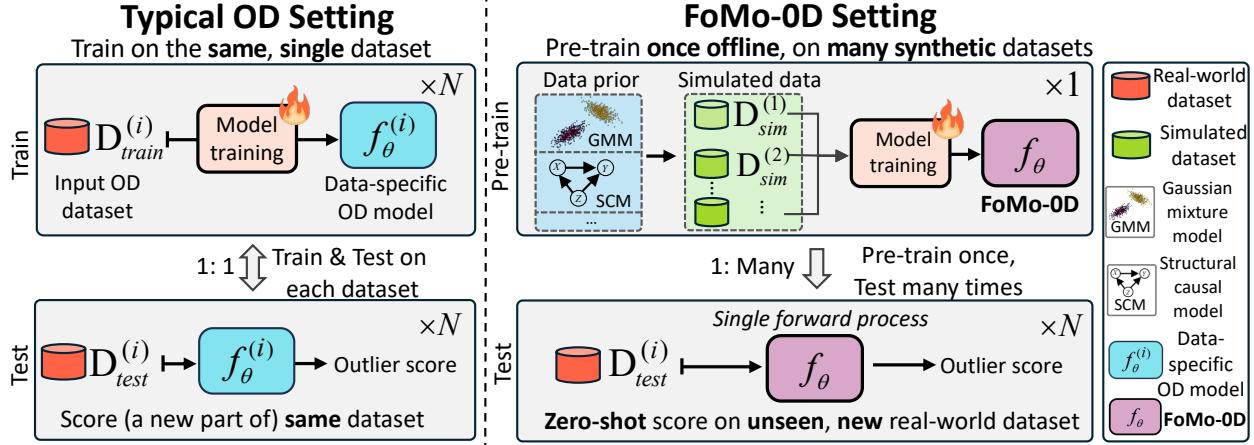


Figure 1: (best in color) Comparison of typical OD vs. the FoMo-0D settings. Given a new unlabeled OD dataset, FoMo-0D not only eliminates the need for model (parameter) training, but most importantly, also abolishes the onerous task of unsupervised model selection (algorithm and hyperparameters).

113 where $p(\phi)$ denotes the prior probability and $p(\mathcal{D}|\phi)$ is the likelihood of the data \mathcal{D} given ϕ .

114 **PFNs and PPD Approximation:** As obtaining the above PPD is generally intractable, Prior-data Fitted
115 Networks (PFNs) are proposed to approximate the PPD (Müller et al., 2022). Unlike traditional machine
116 learning models that are trained directly on observed datasets, PFNs are pre-trained on simulated datasets
117 that are generated according to a prior distribution. Specifically, it contains the pre-training and inference
118 stages described as the following.

119 *Pre-training on synthetic data.* Massive synthetic datasets are generated for the pre-training stage, by first
120 sampling a hypothesis (i.e., the generating mechanism) $\phi \sim p(\phi)$, and then sampling a dataset $\mathcal{D} \sim p(\mathcal{D}|\phi)$.
121 For training, each dataset \mathcal{D} can be split as $\mathcal{D}_{\text{test}} \subset \mathcal{D}$ and $\mathcal{D}_{\text{train}} = \mathcal{D} \setminus \mathcal{D}_{\text{test}}$. Thus, the PFN with parameters
122 θ can be optimized by making predictions on data points in $\mathcal{D}_{\text{test}}$. For a test point $(\mathbf{x}_{\text{test}}, y_{\text{test}}) \in \mathcal{D}_{\text{test}}$, the
123 training loss is formulated as follows.

$$\mathcal{L} = \mathbb{E}_{(\{(\mathbf{x}_{\text{test}}, y_{\text{test}})\} \cup \mathcal{D}_{\text{train}}) \sim p(\mathcal{D})} [-\log q_{\theta}(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})]. \quad (2)$$

124 The above loss can also be interpreted as minimizing the expected KL divergence between $p(\cdot | \mathbf{x}, \mathcal{D})$ and
125 $q_{\theta}(\cdot | \mathbf{x}, \mathcal{D})$ (Müller et al., 2022). In practice, a PFN model q_{θ} is typically implemented by a Transformer-based
126 architecture (Vaswani et al., 2017), which takes $(\mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ as input, where $\mathbf{x}_{\text{test}} \in \mathcal{D}_{\text{test}}$ and $\mathcal{D}_{\text{train}}$ contains
127 an arbitrary number of instances. The output is the conditional class probabilities for \mathbf{x}_{test} . As the whole
128 training set $\mathcal{D}_{\text{train}}$ is passed as input/context to the Transformer, it learns to predict class labels through
129 sample-to-sample attention.

130 *Inference on real-world data.* In the inference stage, a fresh real-world dataset $\mathcal{D}_{\text{train}}$ and some test instance
131 \mathbf{x}_{test} are fed into the (frozen) pre-trained model, which computes the PPD $q_{\theta}(\cdot | \mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ in a single forward
132 pass. Importantly, PFNs do not require gradient-based parameter tuning on new datasets, where prediction
133 is delivered *in less than a second* (Hollmann et al., 2023).

134 In summary, PFNs are trained once, and can be used many times for zero-shot inference on new datasets
135 with different characteristics. The main benefit is that **no training or tuning** is required at the inference
136 stage. This type of learning ability is also termed as in-context learning (ICL) (Xie et al., 2021), which was
137 shown to be effective for various tasks in languages (Brown et al., 2020). In fact, ICL with PFNs is recently
138 shown to be a promising paradigm for supervised classification on tabular datasets (Hollmann et al., 2023).

139 3 FoMo-0D: A Foundation Model for Zero-shot Outlier Detection

140 Inspired by PFNs (Müller et al., 2022; Hollmann et al., 2023; Dooley et al., 2023), we propose **FoMo-0D** for
141 zero-shot OD, which is pre-trained on large-scale synthetic OD datasets for zero-shot detection at inference

¹⁴² time. FoMo-0D eliminates the need for model training on a new dataset and for model selection (both
¹⁴³ algorithm and HPs), which is difficult without any labeled data. The new FoMo-0D paradigm (right) versus
¹⁴⁴ the typical OD setting (left) is illustrated in Figure 1.

¹⁴⁵ In the following we describe our OD data prior, training on prior-simulated datasets, inference on new
¹⁴⁶ datasets, and model architecture and improvements for scalability.

¹⁴⁷ 3.1 Designing a Data Prior for Outlier Detection

¹⁴⁸ Foundation models benefit from massive amounts of datasets available for pre-training, along with high-
¹⁴⁹ capacity model architectures, however, the quantity (and quality) of publicly available tabular OD datasets is
¹⁵⁰ minuscule, compared to the massive size of open-domain text corpora. Even with large quantities of data,
¹⁵¹ Ansari et al. (2024) show that using synthetic data in combination with real-world data improves the overall
¹⁵² zero-shot performance for time-series foundation models. Hence, we design a new data prior from which we
¹⁵³ simulate numerous OD datasets for pre-training FoMo-0D.

¹⁵⁴ Ideally, the data prior should reflect distributions as general and diverse as seen in the real world, however,
¹⁵⁵ “finding a prior supporting a large enough subset of possible [data generating] functions isn’t trivial” (Nagler,
¹⁵⁶ 2023). Surprisingly, our results show that a straightforward, simple-to-implement data prior is sufficient to
¹⁵⁷ achieve remarkable performance.

¹⁵⁸ **Inlier synthesis:** We simulate inliers by drawing from a Gaussian Mixture Model (GMM) with m -clusters
¹⁵⁹ in d -dimensions, with centers $\mu_{jk} \in [-5, 5]$, $j \in [m]$, $k \in [d]$ and diagonal² Σ_j with entries in $(0, 5]$. We create
¹⁶⁰ different GMMs with varying $m \leq M$ and $d \leq D$ chosen uniformly at random from $[M]$ and $[D]$, respectively.
¹⁶¹ From each GMM, we draw a set of S inliers, defined as instances within the 90th percentile of the GMM.

¹⁶² **Outlier synthesis:** Following Han et al. (2022), we generate subspace outliers by first drawing a subset of
¹⁶³ dimensions \mathcal{K} at random, for $|\mathcal{K}| \leq d$, and then generate S points from the “inflated” GMM, which shares
¹⁶⁴ the same centers μ_j ’s with the original GMM but with the inflated (diagonal) covariances $5 \times \Sigma_{j,kk}$ ’s for
¹⁶⁵ $k \in \mathcal{K}$. Outliers are defined as points sampled outside the 90th percentile of the original GMM, which are
¹⁶⁶ labeled based on their Mahalanobis distances (see Property B.6 in the Appendix).

¹⁶⁷ Specifically, we simulate datasets containing $2S = 10,000$ samples (half inlier, half outlier) from the two
¹⁶⁸ corresponding GMMs (original and inflated) with up to $M = 5$ clusters and up to $D = 100$ dimensions.
¹⁶⁹ Example 2-d synthetic datasets are illustrated in Appendix A.

¹⁷⁰ **Remarks:** Our model is not trained on any real-world data but rather, on purely synthetic data (although
¹⁷¹ future work can combine existing benchmark OD datasets with synthesized data, as was done by Ansari
¹⁷² et al. (2024) for time series). While we have intended to extend our preliminary attempt toward designing
¹⁷³ a sophisticated data prior for OD, we found (to our surprise) that even with a basic, GMM-based prior,
¹⁷⁴ FoMo-0D generalizes remarkably well to real-world OD datasets downstream³, outperforming numerous SOTA
¹⁷⁵ baselines. Therefore, we present FoMo-0D with this simple prior to showcase the prowess of PFNs for OD. We
¹⁷⁶ leave as future work the exploration of other priors (Hollmann et al., 2023) and other outlier types (contextual,
¹⁷⁷ dependency, etc. (Steinbuss & Böhm, 2021)), the impact of different priors on performance, as well as prior
¹⁷⁸ mixture composition to further improve performance.

¹⁷⁹ 3.2 (Pre)Training and Inference

¹⁸⁰ **Model (Pre)Training (Once, Offline):** FoMo-0D is a Prior-data Fitted Network (PFN, see Section 2.2)
¹⁸¹ based on the Transformer architecture. In the synthetic prior-data fitting phase, it is trained on datasets
¹⁸² drawn from our OD data prior for tabular data introduced in Section 3.1. Each dataset is simulated from
¹⁸³ a different GMM configuration based on randomly drawn parameters, and consists of varying number of
¹⁸⁴ training samples and dimensions to capture the diversity in real-world tabular datasets. Details are outlined
¹⁸⁵ in Algorithm 1 in Appendix C, and described as follows.

²In early experiments, we found no difference in test performance on synthetic datasets between using diagonal vs. non-diagonal Σ , yet, it is easier to invert diagonal Σ for data synthesis.

³We refer to Appendix G.2 and Figure 16 for an exploration of FoMo-0D performance and GMM goodness of fit on real-world OD datasets.

186 At each time, we first draw a hypothesis (i.e. GMM configuration) uniformly at random, that is, $\phi =$
 187 $\{d \in [D], m \in [M], \{\boldsymbol{\mu}_j\}_{j=1}^m \in [-5, 5]^d, \{\boldsymbol{\Sigma}_j\}_{j=1}^m; \text{diag}(\boldsymbol{\Sigma}_j) \in [-5, 5]^d\}$, and then generate a synthetic
 188 dataset $\mathcal{D} = \{\mathcal{D}_{\text{in}}, \mathcal{D}_{\text{out}}\}$ containing synthetic inlier and outlier samples from the drawn hypothesis and its
 189 variance-inflated variant, respectively.

190 We optimize FoMo-0D’s parameters θ to make predictions on $\mathcal{D}_{\text{test}} = \{\mathcal{D}_{\text{test}}^{\text{in}}, \mathcal{D}_{\text{test}}^{\text{out}}\}$, conditioned on the
 191 inlier-only training data $\mathcal{D}_{\text{train}} \subset \mathcal{D}_{\text{in}}$ based on the cross-entropy loss (see Eq. (2)). During training, $\mathcal{D}_{\text{test}}$
 192 contains a *balanced* number of inlier and outlier samples, where $\mathcal{D}_{\text{test}}^{\text{in}} = \mathcal{D}_{\text{in}} \setminus \mathcal{D}_{\text{train}}$, and $\mathcal{D}_{\text{test}}^{\text{out}} \subset \mathcal{D}_{\text{out}}$ contains
 193 an equal number of samples as $\mathcal{D}_{\text{test}}^{\text{in}}$. To vary the training data size, we subsample $\mathcal{D}_{\text{train}}$ of randomly drawn
 194 size $n \in [n_L, n_U]$, where n_L and n_U denote the lower and upper bounds. In our implementation, we use
 195 $n_L = 500$, and $n_U = 5,000$.

196 FoMo-0D is trained on 200,000 batches (200 epochs \times 1,000 steps/epoch) of $B = 8$ generated datasets in each
 197 batch. While this pre-training phase can be expensive, it is done *only once, offline*. Moreover, we introduce
 198 several scalability improvements to speed up pre-training, as discussed later in Section 3.3. Full details on
 199 the training and implementation of FoMo-0D are given in Appendix C.

200 **Zero-shot Inference (on Unseen/New Dataset):** At inference, the pre-trained FoMo-0D can be employed
 201 on any unseen real-world dataset. Specifically, for a new [unsupervised OD task with inlier-only training data](#)
 202 $\mathcal{D}_{\text{train}}$ and mixed test data $\mathcal{D}_{\text{test}}$, feeding $\langle \mathcal{D}_{\text{train}}, \mathbf{x}_{\text{test}} \rangle$ as input to FoMo-0D (for each $\mathbf{x}_{\text{test}} \in \mathcal{D}_{\text{test}}$ separately)
 203 yields the PPD $q_{\theta}(y|\mathbf{x}_{\text{test}}, \mathcal{D}_{\text{train}})$ in a *single forward pass*. As such, FoMo-0D performs model “training” and
 204 prediction simultaneously at test time. In fact, as the training data is passed as context, FoMo-0D leverages
 205 in-context learning (ICL) (Xie et al., 2021; Garg et al., 2022) for inference.

206 **Remarks:** The key contribution of FoMo-0D goes beyond [eliminating the need for model training for a new](#)
 207 [dataset, it renders model selection an obsolete concern for OD](#). In other words, a practitioner with a new
 208 detection task no longer needs to choose an OD model to train or grapple with tuning any [hyperparameters](#)
 209 [of the said model](#). Further, the speedy, easily parallelizable inference (for *less-than-a-second* per test sample)
 210 is the “icing on the cake”. Figure 1 (right) illustrates (top) pre-train and (bottom) test phases of FoMo-0D,
 211 where the pre-trained FoMo-0D is reused during inference on new datasets directly, unlocking zero-shot OD.

212 3.3 Architecture and Scalability

213 **Architecture and sample-to-sample attention:** Like existing PFNs, FoMo-0D is based on the Transformer
 214 (Vaswani et al., 2017), encoding each sample’s feature vector as a [fixed size token through a linear embedding](#)
 215 [layer \(see Appendix C.2\)](#), and allowing token representations to attend to each other, hence enabling sample-
 216 to-sample attention. We also adopt the three customizations from TabPFN (Hollmann et al., 2023), which
 217 (1) computes self-attention among all the training samples and only cross-attention from test samples to the
 218 training samples, (2) enables varying feature dimensionality by zero-padding, and (3) randomly permutes
 219 input samples while omitting positional encodings to achieve model invariance in the dataset.

220 Given $\mathcal{D}_{\text{train}} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, each self-attention layer outputs n embeddings $\{\mathbf{z}_i\}_{i=1}^n$; where the i -th token is
 221 mapped via linear transformations to a key \mathbf{k}_i , query \mathbf{q}_i and value \mathbf{v}_i , where the i -th output is computed as

$$\mathbf{z}_i = \sum_{j=1}^n \text{softmax}(\{\langle \mathbf{q}_i, \mathbf{k}_{j'} \rangle\}_{j'=1}^n)_j \cdot \mathbf{v}_j . \quad (3)$$

222 The sample-to-sample attention is intriguing from the perspective of OD: many classical OD algorithms
 223 (Aggarwal, 2013) are based on nonparametrics; in particular, they leverage the distances to the k *nearest*
 224 neighbors (k NNs) of a point to compute its outlierness, where k is a critical hyperparameter. One can think
 225 of FoMo-0D as mimicking non-parametric models but by using parametric attention mechanisms. Interestingly,
 226 PFNs are much more robust and flexible than k NN based OD approaches, for (1) sample-to-sample relations
 227 are not pre-specified but rather learned through attention weights, and thus (2) they are not limited to just
 228 the nearest neighbors but rather can *learn which* training points are worth attending to, and (3) as attention
 229 is dataset-wide across all points, there is no need for specifying a cut-off HP value like k , to which most k NN
 230 based OD techniques are sensitive to (Aggarwal & Sathe, 2015; Campos et al., 2016; Goldstein & Uchida,
 231 2016; Ding et al., 2022). We present analyses on sample-to-sample attention in Appendix E.

- To seize the power of scale, we incorporate a scalable architecture and data synthesis into our design to benefit pre-training and inference, as we describe next. The scale-up unlocks a larger context size for FoMo-0D, enabling pre-training and inference on larger datasets with fast speed.
- Scaling up attention with “routers”:** The $\mathcal{O}(n^2)$ quadratic sample complexity at pre-training presents an obstacle for achieving high performance at inference, as it limits pre-training to relatively small training datasets, and degenerates in-context learning that typically benefits from longer context (Xie et al., 2021).
- Toward a high-performance model, we scale up FoMo-0D’s attention via the “router mechanism” of Zhang & Yan (2023). As shown in Figure 2, the main idea is to learn a small number ($R \ll n$) of “routers”, which gather information from all n samples and then distribute the information back to the n output embeddings, in effect, reducing complexity from $\mathcal{O}(n^2)$ to $\mathcal{O}(2Rn) = \mathcal{O}(n)$. This design allows FoMo-0D to **scale linearly** with respect to both dimensionality d and dataset size n in pre-training **as well as during inference**.

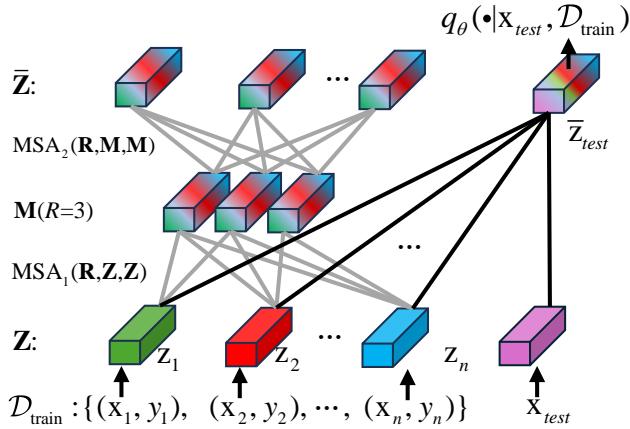


Figure 2: FoMo-0D architecture employs the “router mechanism” for scalable attention.

- Concretely, the representatives first aggregate information from all samples by serving as the query in the multi-head self-attention (MSA):

$$\mathcal{M} = \text{MSA}_1(\mathbf{R}, \mathbf{Z}, \mathbf{Z}), \quad (4)$$

where $\mathbf{R} \in \mathbb{R}^{R \times d}$ depicts the *learnable* vector array of representatives and \mathcal{M} denotes the aggregated messages. Then, the routers distribute the received information among samples by using the sample embeddings as query and the aggregated messages as both key and value:

$$\hat{\mathbf{Z}} = \text{MSA}_2(\mathbf{Z}, \mathcal{M}, \mathcal{M}). \quad (5)$$

Finally, we obtain $\bar{\mathbf{Z}} = \text{LayerNorm}(\hat{\mathbf{Z}} + \mathbf{Z})$ after layer normalization. Note that the test samples only attend to the training samples’ embeddings, computed in the described manner across layers, and are finally fed into the prediction head to estimate the PPD at the output layer.

Scaling up (pre)training data synthesis with linear transforms: Besides the scalability challenge associated with architecture/attention, another computational challenge in pre-training FoMo-0D arises from drawing samples from the data prior, which requires considerable time, especially in high dimensions⁴, provided the large number of datasets we sample (specifically, we utilize a batch size of 8 datasets over 1,000 steps each for 200 epochs).

To give an idea, sampling a dataset with $n = 10,000$ points in $d = 100$ dimensions using 10 CPUs in parallel takes ≈ 0.4 seconds (see Appendix Figure 7). Across 200 training epochs with 1,000 steps each, it adds up to more than 177 hours just to generate 1,6 million datasets on-the-fly. Of course, one can trade storage with compute-time by generating all these datasets apriori via massive parallelism. Nevertheless, synthetic data generation demands considerable time (and/or storage).

⁴This is because the inverse of the $(d \times d)$ covariance matrix plays a crucial role in the process of drawing samples from GMMs, which has $\mathcal{O}(d^3)$ time complexity. (It is also the reason why diagonal Σ_j ’s are favored in our data prior.) In addition, Mahalanobis distance for labeling inliers/outliers also requires the inverse.

262 To scale up data synthesis, FoMo-0D employs two distinct strategies. **First**, we propose *reuse at epoch level*:
 263 that is, one can reuse the same 8K (8×1000) unique datasets at every epoch, or in general, the same $8K \times P$
 264 datasets periodically at every P epochs. A larger P would lead to more diversity in terms of the overall
 265 pre-training data used.

266 **Second**, we propose *reuse at dataset level via transformation*: that is, having generated one unique dataset
 267 $\mathbf{X} \in \mathbb{R}^{n \times d}$ from a GMM, we propose a linear transform $T(\mathbf{x})$ of the form $\mathbf{Wx} + \mathbf{b}$ for randomly drawn
 268 parameters $\mathbf{W} \in \mathbb{R}^{d \times d}$ and $\mathbf{b} \in \mathbb{R}^d$ (see Appendix B.1).⁵ This simple yet efficient transformation creates a
 269 new dataset, akin to one being drawn from another GMM with centers $T(\boldsymbol{\mu}_j) = \mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}$ and covariance
 270 $T(\boldsymbol{\Sigma}_j) = \mathbf{W}\boldsymbol{\Sigma}_j\mathbf{W}^T, \forall j \in [m]$. Note that we do not actually materialize these parameters but only transform
 271 the dataset. As we show in the following, such transformations preserve the Mahalanobis distances as well as
 272 the percentile thresholds for labeling points as inlier/outlier. Details and proofs are given in Appendix B.

273 **Lemma 3.1.** *Linear transform T with invertible \mathbf{W} on \mathcal{G}_m^d preserves Mahalanobis distances.*

274 **Lemma 3.2.** *Linear transform T with invertible \mathbf{W} on \mathcal{G}_m^d preserves the percentiles of the GMM.*

275 The implication of these lemmas is that a linear transformation of a dataset from a GMM retains the identity
 276 of the inliers and outliers, i.e. no relabeling is required. Moreover, notice that as a byproduct we obtain
 277 a transformed dataset as though it is drawn from a GMM with a *non-diagonal* covariance matrix which,
 besides the time savings, offers a slightly more complex data prior.

278 To reach 8K unique datasets for each epoch, we first generate 500 datasets from different GMMs (with
 279 varying configurations), then employ 15 different linear transformations to each dataset by varying \mathbf{W} and \mathbf{b} .
 280 Drawing each (\mathbf{W}, \mathbf{b}) takes ≈ 0.02 seconds, while the matrix-matrix product of \mathbf{X} ($n \times d$) and \mathbf{W} ($d \times d$)
 281 takes negligible time (for $d \leq 100$). Thus, obtaining a transformed dataset offers $20 \times$ speed-up compared to
 282 generating one (0.02 vs. 0.4 seconds).

283 4 Experiments

284 4.1 Setup

285 We present the experiment setup briefly, including data synthesis, real-world datasets, baselines, metrics and
 286 HPs. For more details, we refer to Appendix D.

287 **Pre-training Dataset Synthesis:** During pre-training, we generate unique GMM datasets by first drawing
 288 a configuration, including dimensionality $d \in [D]$, number of components $m \in [M]$, centers $\{\boldsymbol{\mu}_j\}_{j=1}^m$ (each
 289 $\boldsymbol{\mu}_j \in [-5, 5]^d$) and covariances $\{\boldsymbol{\Sigma}_j\}_{j=1}^m$ ($\text{diag}(\boldsymbol{\Sigma}_j) \in [-5, 5]^d$). We set $M = 5$ and vary $D \in \{20, 100\}$ to
 290 study pre-training with relatively small and high dimensional datasets, respectively. We synthesize inliers
 291 and outliers described in Section 3.1.

292 **Real-world Benchmark Datasets:** While pre-training is purely on synthetic datasets, we evaluate
 293 FoMo-0D on **57** real-world datasets from ADBench (Han et al., 2022) (see Table 20 in Appendix J). Following
 294 Livernoche et al. (2024), we use 5 train/test splits of each dataset via different seeds and report mean
 295 performance and standard deviation. Note that the baselines require model re-training and inference for each
 296 $\mathcal{D}_{\text{train}}/\mathcal{D}_{\text{test}}$ split, while FoMo-0D uses the splits only for inference as $\mathcal{D}_{\text{train}}$ is passed as context.

297 **Baselines:** We compare FoMo-0D against **26** baselines, from classical/shallow methods to modern/deep
 298 models. The baselines are imported from one of the latest papers that proposed the SOTA diffusion-based
 299 OD model, DTE (Livernoche et al., 2024), and three variants; DTE-C, DTE-IG, DTE-NP. As such, the long
 300 list of baselines we compare to constitutes one of the most comprehensive in the literature. We refer to the
 301 original paper for more details.

302 **Model Implementation:** We train our final model for 200,000 steps with a batch size of 8 datasets. That is,
 303 FoMo-0D is trained on 1,600,000 synthetically generated datasets. This training takes about 25 hours on 1 GPU
 304 (Nvidia RTX A6000). Each dataset had a fixed size of 10,000 samples, with $|\mathcal{D}_{\text{train}}| \in [n_L = 500, n_U = 5000]$,

⁵In practice, we apply the linear transform on the subspace of inflated features only, wherein inliers and outliers are defined, which remains to be a multi-variate GMM.

305 and the rest as $\mathcal{D}_{\text{test}}$ with balanced number of inliers and outliers. Other details of FoMo-0D, including the
 306 training algorithm, model architecture, data synthesis and reuse, and hardware are in Appendix C.

307 **Metrics and Hypothesis Testing:** Detection performance is w.r.t. 3 widely-used metrics for OD: AUROC;
 308 area under ROC curve, AUPR; area under Precision-Recall curve, and F1 score; using threshold at the true
 309 number of outliers in the test data (varies by dataset) Livernoche et al. (2024).

310 To compare different methods on ADBench, we compute their rank on each dataset (lower is better), and
 311 present average rank across datasets. This is an alternative to the average metric (e.g. AUROC), which is not
 312 meaningful when tasks vary widely in terms of their difficulties.

313 In addition, we perform significance tests to compare two methods statistically, using the one-sided paired
 314 Wilcoxon signed rank test (Demšar, 2006) between FoMo-0D and a baseline based on the performances across
 315 all datasets, with the alternative hypothesis suggesting the “baseline-minus-FoMo-0D” performance gap is
 316 greater than zero. We consider results to be significant at $\alpha = 0.05$ following convention.

317 **Hyperparameters (HPs):** Importantly, Livernoche et al. (2024) picked for each baseline the best-
 318 performing set of HPs as recommended by the authors in their original paper. As for their own DTE, which
 319 behaves similarly to kNN, they use $k = 5$ and set the same k for the kNN baseline (Ramaswamy et al., 2000)
 320 to be consistent. However, it is well known that kNN is sensitive to the value of k (Aggarwal & Sathé, 2015),
 321 and so are many other OD models to their respective HPs (Campos et al., 2016; Goldstein & Uchida, 2016;
 322 Zhao et al., 2021; Ding et al., 2022).

323 Therefore, besides comparing FoMo-0D with the 26 baselines in Livernoche et al. (2024), respectively for
 324 AUROC, F1, and AUPR (Livernoche et al., 2024), we also compare to the top-4⁶ best-performing baselines
 325 (in order: DTE-NP, kNN, ICL, and DTE-C) on their average performance across a list of different HP settings.
 326 Such an approach reflects their expected performance under HP values selected at random, in the absence of
 327 any other prior knowledge, as recommended by Goldstein & Uchida (2016) “to get a fair evaluation when
 328 comparing [OD] algorithms”. We annotate the method name with ^{avg} for the version with performance
 329 averaged over varying HPs. The detailed list of HP values for each top baseline is given in Appendix D.4.
 330 Overall, we compare FoMo-0D to 30 baselines; 26 from Livernoche et al. (2024) and ^{avg} of the top-4.

331 4.2 Results

332 **Detection performance:** Table 1 presented the comparison of FoMo-0D w/ $D = 100$ to all baselines
 333 w.r.t. average rank across datasets as well as pairwise Wilcoxon signed rank tests based on AUROC, and
 334 we present full results on all datasets and all metrics in Appendix I. We find that among 30
 335 baselines and 2 variants of FoMo-0D (w/ $D = 100$ and $D = 20$), FoMo-0D w/ $D = 100$ performs as well as the
 336 2nd best model (kNN with default HP; $k = 5$) on all datasets. While DTE-NP outperforms FoMo-0D with
 337 author-recommended $k = 5$, we find that DTE-NP^{avg} is on par with FoMo-0D.

338 In our tests, $p > \alpha = 0.05$ implies no statistical evidence for performance difference between two methods.
 339 FoMo-0D w/ $D = 100$ performs statistically no different from all baselines on datasets with $d \leq 100$ (i.e., “at
 340 its own game” when pre-training data dimensions align with real-world datasets), while it outperforms the
 341 majority of baselines (where $p > 1 - \alpha$). These results continue to hold on datasets with $d \leq 500$.

342 Table 2 shows similar results for FoMo-0D w/ $D = 20$, which is pre-trained on datasets with considerably
 343 fewer dimensions. Even in this limited setting, it performs on par with the 3rd best baseline (ICL, with
 344 default HP) against 30 baselines, with an increased p -value (0.437) when compared to ICL^{avg}. On datasets
 345 with $d \leq 20$ which align with its pre-training data, it outperforms the top 5th baseline and the majority of
 346 others. With FoMo-0D pre-trained purely on synthetic datasets from a simple prior in small dimensions, these
 347 results showcase the prowess of PFNs for OD.

348 Figure 3 shows the distribution of AUCROC across datasets for all models (See the distribution of ranks with
 349 respect to AUCROC in Appendix Figure 17). While p-values are the most statistically conclusive, FoMo-0D
 350 achieves a relatively small average rank with notably higher AUCROC across datasets compared to the

⁶ To rank the 26 baselines, we compute the 26×26 p -values of the pairwise Wilcoxon signed rank test (see Appendix Figure 24), and order them by their mean p -value against other baselines.

Table 2: p -values of the one-sided Wilcoxon signed rank test, comparing FoMo-0D (w/ $D = 20$) to **top 10** baselines with default HPs, and **top 4^{avg}** baselines⁶ with **avg.** performance over varying HPs (denoted w/ ^{avg}) over All (57) datasets, those (24) w/ $d \leq 20$ and (38) datasets w/ $d \leq 50$ dimensions. Although pretrained on datasets w/ small $D = 20$, FoMo-0D shows **no statistically significant difference from the top 3rd baseline** (ICL, w/ $p = 0.089$) over All datasets, while it outperforms (w/ $p > 1 - \alpha$) the top 5th (LOF) and onward baselines over datasets w/ $d \leq 20$ (aligned w/ pretraining where $D = 20$) and on datasets w/ $d \leq 50$ (generalizing beyond pretraining). We use underline and **bold** to indicate $p < \alpha$ and $p > 1 - \alpha$. Rank is avg.'ed over all 57 datasets, where methods are ranked on each dataset w.r.t. AUROC. (experiment setting: $D = 20$, $P = 50$, $R = 500$, train/inference context size=5K, no data transformation, $\alpha = 0.05$)

	FoMo-0D	DTE-NP	<i>k</i> NN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP ^{avg}	<i>k</i> NN ^{avg}	ICL ^{avg}	DTE-C ^{avg}
$d \leq 20$	-	0.572	0.789	0.968	0.616	0.993	0.989	1.000	0.978	0.906	0.992	0.813	0.924	0.999	1.000
$d \leq 50$	-	0.347	0.794	0.893	0.946	0.997	0.988	1.000	0.963	0.994	0.986	0.574	0.847	0.995	1.000
All	-	<u>0.001</u>	<u>0.019</u>	0.089	0.159	0.394	0.434	0.703	0.516	0.752	0.679	<u>0.007</u>	0.062	0.437	1.000
Rank(avg)	12.59	7.19	8.57	10.34	10.79	11.82	12.81	12.8	12.52	13.50	13.34	8.60	10.63	12.44	21.43

351 majority of the baselines. Appendix H presents another comparison between detectors through performance
 352 profile plots (Dolan & Moré, 2002).

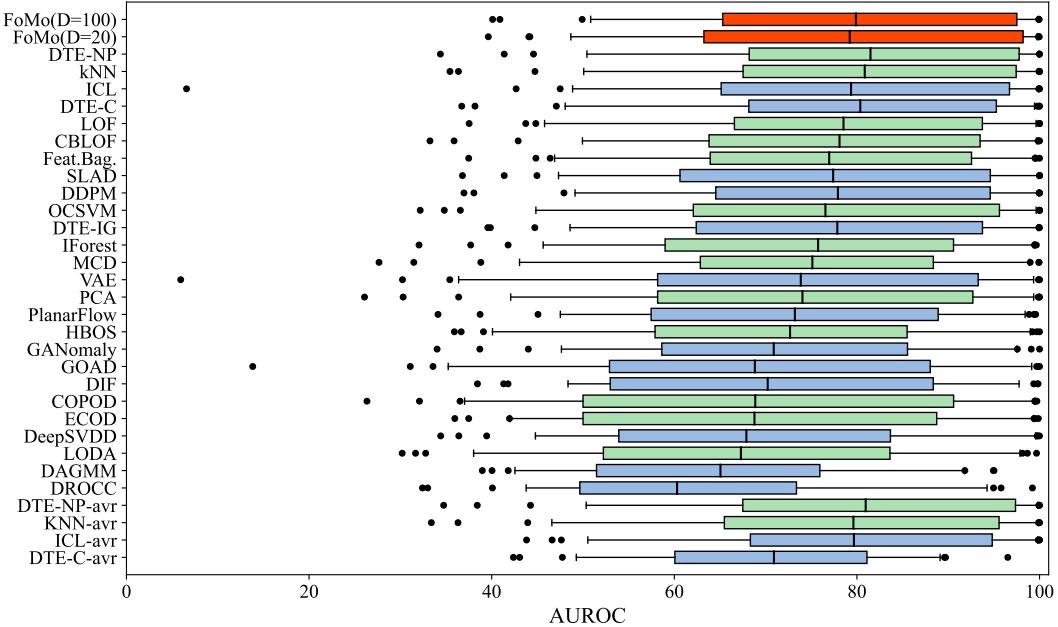


Figure 3: (best in color) AUROC (higher is better) distribution across all **57** real-world datasets shown via boxplots for (from top to bottom) FoMo-0D in red, all **26** baselines ordered by mean p -value⁶ (shallow and deep baselines in green and blue), and **top 4** baselines' ^{avg} variants. The vertical line depicts the mean, the box shows the 25-75%, bars range 5-95%, and circles show the datasets at the tails.

353 **Running time:** Table 3 presents the total training time and the average inference time per test sample
 354 as measured on the largest dataset for FoMo-0D and the top-3 baselines. Given a new dataset, FoMo-0D
 355 bypasses model training (and HP tuning) and directly performs inference, with an average of 7.7 ms per
 356 sample (see Appendix Figure 6). In comparison, all baseline methods need to train on each individual dataset
 357 preceding inference. This training time can be high for deep learning based models like ICL, and further
 358 compounded with training multiple models for hyperparameter tuning purposes. Even for non-parametric
 359 and/or shallow models like *k*NN and DTE-NP (which queries *k* nearest neighbors), the training involves
 360 various data pre-processing steps such as constructing a tree-like data structure for fast (often approximate)
 361 *k*NN distance querying.

Table 3: Train-time and Inference-time (in milliseconds) of **FoMo-0D** and the top-3⁶ baselines (w/ *default* HPs, *excluding* the time for model selection/hyperparameter optimization) on our largest dataset **donors** (see Appendix Table 20). **FoMo-0D** skips any model training or fine-tuning and takes a mere forward pass for inference out-of-the-box.

Method	FoMo-0D	DTE-NP	kNN	ICL
Train-t (total)	none/0-shot	56.83	1433.74	186461.48
Infer-t (per sample)	7.7	0.76	0.17	0.01

362 4.3 Ablation Analyses

363 Extensive ablations in Appendix F analyze the effect of D in F.1, cost and performance by varying R in F.2
 364 and F.3, context size in F.4, reuse periodicity P in F.5, effect of data transformation T on performance and
 365 speed-up in F.6 and F.7, data diversity and prolonged training in F.8, quantile transformation in F.9, and
 366 **model size in F.10**.

367 4.4 Generalization Analyses

368 Our results that synthetic pretraining at large, and GMM data prior in particular, enables **FoMo-0D** to reach
 369 remarkable performance on real world tasks call for deeper investigation on its generalization. To this end,
 370 Appendix G provides an extensive analysis on **FoMo-0D**'s generalization to out-of-distribution (OOD) synthetic
 371 datasets in G.1, GMM statistical goodness-of-fit analyses of real-world datasets in ADBench in G.2, as well
 372 as generalization to real-world OOD detection tasks in G.3.

373 5 Related Work

374 **Outlier Detection (OD):** Thanks to diverse applications in numerous fields, such as security, finance,
 375 manufacturing, to name a few, OD on tabular (or point-cloud) datasets has a vast literature with a long list
 376 of techniques. For earlier, shallow approaches preceding the advances in deep learning, we refer to the books
 377 by Aggarwal (2013) and Aggarwal & Sathe (2017). The modern, deep learning based techniques are surveyed
 378 in Chalapathy & Chawla (2019); Pang et al. (2021); Ruff et al. (2021). Most recent deep OD techniques take
 379 advantage of newly emerging paradigms, including self-supervised learning (Hojjati et al., 2022; Yoo et al.,
 380 2023) as well as the most recently popularized diffusion-based models (Yoon et al., 2023; Livernoche et al.,
 381 2024; Du et al., 2024; He et al., 2024).

382 **Unsupervised Model Selection for OD:** It is typical of models to exhibit various hyperparameters (HPs)
 383 that play a role in the bias-variance trade-off and hence the generalization performance, and OD models are
 384 no exception. Many earlier work on OD showed the sensitivity of classical (i.e. shallow) OD methods to the
 385 choice of their HP(s) (Aggarwal & Sathe, 2015; Campos et al., 2016; Goldstein & Uchida, 2016). Similarly,
 386 sensitivity to HPs has also been shown for deep OD models more recently (Zhao et al., 2021; Ding et al.,
 387 2022), as well as for those relying on self-supervised learning/data augmentation (Yoo et al., 2023).

388 While critical, work on unsupervised outlier model selection (UOMS) is slim as compared to the vast literature
 389 on detection methods. A handful of existing, mostly heuristic strategies has been studied by Ma et al. (2023)
 390 reporting discouraging results; they have shown that existing heuristics are either not significantly different
 391 from random selection, or do not outperform iForest (Liu et al., 2008) with its default HPs.

392 More recent UOMS approaches go beyond heuristic measures and instead design scalable hyperensembles
 393 (Ding et al., 2022; 2024), as well as take advantage of meta-learning on historical real-world OD datasets
 394 (Zhao et al., 2021; 2022; Zhao & Akoglu, 2024). These approaches demonstrate the value of learning from
 395 many other OD datasets, and transfer these learnings to a new dataset. While sharing the same spirit on
 396 learning from a large collection of (in our case, simulated) datasets, our **FoMo-0D** differs from these prior art
 397 in a key aspect: **FoMo-0D** is *not* a model selection technique, but rather, a foundation model that abolishes
 398 model training and selection altogether. As such, it unlocks zero(0)-shot inference on a new task.

399 **Prior-data Fitted Networks:** Based on the seminal work by Müller et al. (2022), Prior-data-fitted
400 Networks (PFNs) establish a new paradigm for machine learning, where a PFN is pretrained on synthetic
401 datasets generated from a data prior, and the pretrained PFN can then infer the posterior predictive
402 distribution (PPD) for test points in a new dataset in a single forward pass, through in-context learning (Xie
403 et al., 2021; Garg et al., 2022). It is shown that PFNs provably approximate Bayesian inference (Müller
404 et al., 2022). Follow-up TabPFN (Hollmann et al., 2023) and its v2 Hollmann et al. (2025) achieved SOTA
405 classification performance on small tabular datasets of size up to 1024. Other subsequent works designed
406 LC-PFN (Adriaensen et al., 2024) and ForecastPFN (Dooley et al., 2023), respectively zero-shot learning
407 curve extrapolation and zero-shot time-series forecasting models, trained purely on synthetic data. PFN4BO
408 (Müller et al., 2023) employed PFNs for Bayesian optimization, while Nagler (2023) studied the statistical
409 foundations of PFNs. Others proposed scaling the context size to enable training on larger datasets toward
410 better generalization (Ma et al., 2024; Feuer et al., 2023; 2024; Qu et al., 2025).

411 Our proposed FoMo-OD differs from these in being the first PFN for OD, using a novel inlier/outlier data
412 prior, employing linear transform for fast data synthesis, and incorporating the “router” attention mechanism
413 for linear-time scalability w.r.t. context size. See Appendix K for additional details.

414 **Zero-Shot Outlier Detection:** Foundation models pretrained on massive text and image corpora, such as
415 large language and/or vision models (L(V)LMs) like OpenAI’s GPT-series (Achiam et al., 2023), DALL-E
416 (Ramesh et al., 2021) and Flamingo (Alayrac et al., 2022), CLIP (Radford et al., 2021), and LLaVA (Liu
417 et al., 2024) to name a few, have demonstrated remarkable success on several zero-shot tasks in CV and NLP.
418 Follow-up work extended these models for zero-shot out-of-distribution detection (Esmaeilpour et al., 2022),
419 zero-shot image OD (Liznerski et al., 2022; Jeong et al., 2023; Zhou et al., 2024) as well as dialogue-based
420 industrial image anomaly detection (Gu et al., 2024).

421 Foundation models, however, do not exist for tabular data which is widespread across OD applications in
422 the real world, such as detecting credit card fraud, network intrusion, medical anomalies, and any sensor
423 measurement abnormalities, to name a few. The recent ACR model by Li et al. (2023) on zero-shot OD
424 does *not* rely on a pretrained foundation model, but rather is meta-trained on each specific domain using
425 inlier-only datasets from the *same domain*. Concurrent to our work, Li et al. (2024) apply pretrained LLMs
426 for prompt-based OD on tabular data which they serialize to text. Similar to our work, they also use *simulated*
427 labeled OD datasets to fine-tune several existing LLMs to improve their performance. Their work, however,
428 is quite preliminary in several fronts; a key limitation is that they assume independent features and query the
429 LLM one-feature-at-a-time to reach an outlier score. Further, they fine-tune using only 5,000 data batches
430 with up to 100 samples each, subsample 150 points and the first 10 columns of each dataset for evaluation
431 (due to GPU memory constraint), and their testbed includes only two baseline methods. In contrast, FoMo-OD
432 employs and pretrains PFNs at a much larger scale with rigorous evaluation on a much larger testbed.

433 6 Conclusion

434 This work introduced FoMo-OD, **the first foundation model for outlier detection (OD)** on tabular
435 data. It capitalizes on the in-context learning of a Transformer model pre-trained on a large number of
436 synthetic datasets that can then perform **zero-shot** inference on a new dataset, without *any* hyper/parameter
437 tuning/training. FoMo-OD breaks new ground by fully abolishing the notoriously-hard model selection task for
438 unsupervised OD (see Impact Statement). Further, FoMo-OD offers extremely fast inference thanks to a mere
439 single forward pass. Against a long list of **26** SOTA baselines on **57** public real-world datasets, FoMo-OD
440 performs on par with the *2nd* best baseline, while outperforming the majority of the baselines. Future work
441 could expand our data prior and explore similar directions for zero-shot OD beyond tabular data. For a
442 detailed discussion on limitations and future directions, we refer to Appendix L.

443 Broader Impact Statement

444 FoMo-0D offers zero-shot outlier detection (OD), abolishing not only parameter training but also model
 445 selection given a new dataset. This is a radical paradigm shift for the OD literature, which historically focused
 446 on designing new models and recently unsupervised model selection. Obviating the need for either, we expect
 447 FoMo-0D to route attention of the community from new model design and selection to designing better data
 448 priors and gathering datasets for pre-training, along with better and more scalable architectures for PFN.
 449 From the applied perspective, a zero-shot OD model like FoMo-0D is a game changer for practitioners! Given
 450 the plethora of OD algorithms to choose from, each with a list of hyperparameters to set, not having the tools
 451 for effective and efficient model selection leaves the practitioners with a “choice paralysis”. With FoMo-0D,
 452 practitioners can not only bypass such dilemmas on one dataset, but thanks to the “train once, use many times”
 453 nature of pre-trained models, they can do so for any dataset, including those arriving over time. FoMo-0D is
 454 lightweight and optimized for fast inference (without any additional training or tuning), making it especially
 455 attractive for real-time or resource-constrained applications. While attractive from a performance viewpoint,
 456 we remark that FoMo-0D currently does not factor into account metrics beyond detection performance, such
 457 as fairness, biases, or other potential blindspots. In sensitive real-world applications, social and ethical costs
 458 of incorrect detections should be taken into consideration.

459 References

- 460 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo
 461 Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint*
 462 *arXiv:2303.08774*, 2023.
- 463 Steven Adriaensen, Herilalaina Rakotoarison, Samuel Müller, and Frank Hutter. Efficient bayesian learning
 464 curve extrapolation using prior-data fitted networks. *Advances in Neural Information Processing Systems*,
 465 36, 2024.
- 466 Charu C. Aggarwal. *Outlier Analysis*. Springer, 2013.
- 467 Charu C. Aggarwal and Saket Sathe. Theoretical foundations and algorithms for outlier ensembles. *Acm sigkdd explorations newsletter*, 17(1):24–47, 2015.
- 469 Charu C. Aggarwal and Saket Sathe. *Outlier Ensembles - An Introduction*. Springer, 2017.
- 470 Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. Gandomaly: Semi-supervised anomaly detection
 471 via adversarial training. In *ACCV*, pp. 622–637. Springer, 2019.
- 472 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc,
 473 Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for
 474 few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- 475 Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr
 476 Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al. Chronos: Learning
 477 the language of time series. *arXiv preprint arXiv:2403.07815*, 2024.
- 478 Lion Bergman and Yedid Hoshen. Classification-based anomaly detection for general data. In *International
 479 Conference on Learning Representations*, 2020.
- 480 Markus M. Breunig, Hans-Peter Kriegel, Raymond T. Ng, and Jörg Sander. Lof: Identifying density-based
 481 local outliers. In *International Conference on Management of Data*, 2000.
- 482 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
 483 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen
 484 Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris
 485 Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher
 486 Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot
 487 learners. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901, 2020.

- 488 Guilherme O Campos, Arthur Zimek, Jörg Sander, Ricardo JGB Campello, Barbora Micenková, Erich
 489 Schubert, Ira Assent, and Michael E Houle. On the evaluation of unsupervised outlier detection: measures,
 490 datasets, and an empirical study. *Data mining and knowledge discovery*, 30:891–927, 2016.
- 491 George Casella and Roger Berger. *Statistical inference*. CRC Press, 2024.
- 492 Raghavendra Chalapathy and Sanjay Chawla. Deep learning for anomaly detection: A survey. *arXiv preprint*
 493 *arXiv:1901.03407*, 2019.
- 494 Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine learning*
 495 *research*, 7:1–30, 2006.
- 496 Xueying Ding, Lingxiao Zhao, and Leman Akoglu. Hyperparameter sensitivity in deep outlier detection:
 497 Analysis and a scalable hyper-ensemble solution. *Advances in Neural Information Processing Systems*, 35:
 498 9603–9616, 2022.
- 499 Xueying Ding, Yue Zhao, and Leman Akoglu. Fast unsupervised deep outlier model selection with hypernet-
 500 works. *ACM SIGKDD*, 2024.
- 501 Elizabeth D Dolan and Jorge J Moré. Benchmarking optimization software with performance profiles.
 502 *Mathematical programming*, 91:201–213, 2002.
- 503 Samuel Dooley, Gurnoor Singh Khurana, Chirag Mohapatra, Siddartha Venkat Naidu, and Colin White. Fore-
 504 castPFN: Synthetically-trained zero-shot forecasting. In *Thirty-seventh Conference on Neural Information*
 505 *Processing Systems*, 2023.
- 506 Xuefeng Du, Yiyou Sun, Jerry Zhu, and Yixuan Li. Dream the impossible: Outlier imagination with diffusion
 507 models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 508 Sepideh Esmaeilpour, Bing Liu, Eric Robertson, and Lei Shu. Zero-shot out-of-distribution detection based
 509 on the pre-trained model CLIP. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36,
 510 pp. 6568–6576, 2022.
- 511 Benjamin Feuer, Niv Cohen, and Chinmay Hegde. Scaling tabPFN: Sketching and feature selection for
 512 tabular prior-data fitted networks. In *NeurIPS 2023 Second Table Representation Learning Workshop*,
 513 2023.
- 514 Benjamin Feuer, Robin Tibor Schirrmeyer, Valeria Cherepanova, Chinmay Hegde, Frank Hutter, Micah
 515 Goldblum, Niv Cohen, and Colin White. Tunetables: Context optimization for scalable prior-data fitted
 516 networks, 2024.
- 517 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context?
 518 a case study of simple function classes. *Advances in Neural Information Processing Systems*, 2022.
- 519 Markus Goldstein and Andreas Dengel. Histogram-based outlier score (hbos): A fast unsupervised anomaly
 520 detection algorithm. *KI-2012: poster and demo track*, 1:59–63, 2012.
- 521 Markus Goldstein and Seiichi Uchida. A comparative evaluation of unsupervised anomaly detection algorithms
 522 for multivariate data. *PloS one*, 11(4), 2016.
- 523 Sachin Goyal, Aditi Raghunathan, Moksh Jain, Harsha Simhadri, and Prateek Jain. Drocc: Deep robust
 524 one-class classification. In *Proceedings of the 37th International Conference on Machine Learning*, ICML’20.
 525 JMLR.org, 2020.
- 526 Zhaopeng Gu, Bingke Zhu, Guibo Zhu, Yingying Chen, Ming Tang, and Jinqiao Wang. AnomalyGPT:
 527 Detecting industrial anomalies using large vision-language models. In *Proceedings of the AAAI Conference*
 528 *on Artificial Intelligence*, volume 38, pp. 1932–1940, 2024.
- 529 A. Gut. *An Intermediate Course in Probability*. Springer Texts in Statistics. Springer New York, 2009. ISBN
 530 9781441901620.

- 531 Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. Adbench: Anomaly detection
 532 benchmark. *Advances in Neural Information Processing Systems*, 35, 2022.
- 533 Haoyang He, Jiangning Zhang, Hongxu Chen, Xuhai Chen, Zhishan Li, Xu Chen, Yabiao Wang, Chengjie
 534 Wang, and Lei Xie. A diffusion-based framework for multi-class anomaly detection. In *Proceedings of the*
 535 *AAAI Conference on Artificial Intelligence*, volume 38, pp. 8472–8480, 2024.
- 536 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In
 537 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- 538 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural*
 539 *Information Processing Systems*, 33:6840–6851, 2020.
- 540 Hadi Hojjati, Thi Kieu Khanh Ho, and Narges Armanfard. Self-supervised anomaly detection: A survey and
 541 outlook. *arXiv preprint arXiv:2205.05173*, 2022.
- 542 Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. TabPFN: A transformer that
 543 solves small tabular classification problems in a second. In *The Eleventh International Conference on*
 544 *Learning Representations*, 2023.
- 545 Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Ti-
 546 bbor Schirrmeister, and Frank Hutter. Accurate predictions on small data with a tabular foundation model.
 547 *Nature*, 637(8045):319–326, 2025.
- 548 Catherine Huber-Carol, Narayanaswamy Balakrishnan, Mikhail Nikulin, and Mounir Mesbah. *Goodness-of-fit*
 549 *tests and model validity*. Springer Science & Business Media, 2012.
- 550 Jongheon Jeong, Yang Zou, Taewan Kim, Dongqing Zhang, Avinash Ravichandran, and Onkar Dabeer.
 551 Winclip: Zero-/few-shot anomaly classification and segmentation. In *Proceedings of the IEEE/CVF*
 552 *Conference on Computer Vision and Pattern Recognition*, pp. 19606–19616, 2023.
- 553 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray,
 554 Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint*
 555 *arXiv:2001.08361*, 2020.
- 556 Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 557 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- 558 Aodong Li, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt. Zero-shot anomaly
 559 detection via batch normalization. In *Thirty-seventh Conference on Neural Information Processing Systems*,
 560 2023.
- 561 Aodong Li, Yunhan Zhao, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt.
 562 Anomaly detection of tabular data using LLMs. *arXiv preprint arXiv:2406.16308*, 2024.
- 563 Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 Eighth IEEE International*
 564 *Conference on Data Mining*, pp. 413–422, 2008.
- 565 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural*
 566 *information processing systems*, 36, 2024.
- 567 Victor Livernoche, Vineet Jain, Yashar Hezaveh, and Siamak Ravanbakhsh. On diffusion modeling for
 568 anomaly detection. In *ICLR*, 2024.
- 569 Philipp Liznerski, Lukas Ruff, Robert A Vandermeulen, Billy Joe Franks, Klaus-Robert Müller, and Marius
 570 Kloft. Exposing outlier exposure: What can be learned from few, one, and zero outlier images. *arXiv*
 571 *preprint arXiv:2205.11474*, 2022.
- 572 Junwei Ma, Valentin Thomas, Guangwei Yu, and Anthony L. Caterini. In-context data distillation with
 573 tabpfn. *CoRR*, abs/2402.06971, 2024.

- 574 Martin Q Ma, Yue Zhao, Xiaorong Zhang, and Leman Akoglu. The need for unsupervised outlier model
 575 selection: A review and evaluation of internal evaluation strategies. *ACM SIGKDD Explorations Newsletter*,
 576 25(1):19–35, 2023.
- 577 Samuel Müller, Matthias Feurer, Noah Hollmann, and Frank Hutter. PFNs4BO: in-context learning for
 578 bayesian optimization. In *International Conference on Machine Learning*, 2023.
- 579 Samuel Müller, Noah Hollmann, Sebastian Pineda-Arango, Josif Grabocka, and Frank Hutter. Transformers
 580 can do bayesian inference. *ICLR*, 2022.
- 581 Thomas Nagler. Statistical foundations of prior-data fitted networks. In *ICML*, volume 202 of *Proceedings of
 582 Machine Learning Research*. PMLR, 2023.
- 583 Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media,
 584 2012.
- 585 Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. Deep learning for anomaly
 586 detection: A review. *ACM computing surveys (CSUR)*, 54(2):1–38, 2021.
- 587 Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding
 588 important examples early in training. *Advances in neural information processing systems*, 34:20596–20607,
 589 2021.
- 590 Judea Pearl. *Causality*. Cambridge University Press, 2009.
- 591 Jingang Qu, David Holzmüller, Gaël Varoquaux, and Marine Le Morvan. Tabicl: A tabular foundation model
 592 for in-context learning on large data. *arXiv preprint arXiv:2502.05564*, 2025.
- 593 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
 594 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural
 595 language supervision. In *International conference on machine learning*, 2021.
- 596 Sridhar Ramaswamy, Rajeev Rastogi, and Kyuseok Shim. Efficient algorithms for mining outliers from large
 597 data sets. In *ACM SIGMOD international conference on management of data*, 2000.
- 598 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and
 599 Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pp.
 600 8821–8831. Pmlr, 2021.
- 601 Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the emergence
 602 of non-bayesian in-context learning for regression. *Advances in Neural Information Processing Systems*,
 603 2024.
- 604 Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *Proceedings of
 605 the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning
 606 Research*, pp. 1530–1538, Lille, France, 07–09 Jul 2015. PMLR.
- 607 Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder,
 608 Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *International Conference on Machine
 609 Learning*, pp. 4393–4402, 2018.
- 610 Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius
 611 Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly
 612 detection. *Proceedings of the IEEE*, 109(5):756–795, 2021.
- 613 Tom Shenkar and Lior Wolf. Anomaly detection for tabular data with internal contrastive learning. In
 614 *International Conference on Learning Representations*, 2022.
- 615 Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling
 616 laws: beating power law scaling via data pruning. *Advances in Neural Information Processing Systems*, 35:
 617 19523–19536, 2022.

- 618 Georg Steinbuss and Klemens Böhm. Benchmarking unsupervised outlier detection with realistic synthetic
 619 data. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(4):1–20, 2021.
- 620 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix,
 621 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard
 622 Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. URL
 623 <https://arxiv.org/abs/2302.13971>.
- 624 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser,
 625 and Illia Polosukhin. Attention is all you need. In *NIPS*, pp. 5998–6008, 2017.
- 626 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning
 627 as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- 628 Hongzuo Xu, Yijie Wang, Juhui Wei, Songlei Jian, Yizhou Li, and Ning Liu. Fascinating supervisory signals
 629 and where to find them: Deep anomaly detection with scale learning. In *International Conference on
 630 Machine Learning*, pp. 38655–38673. PMLR, 2023.
- 631 Jaemin Yoo, Tiancheng Zhao, and Leman Akoglu. Data augmentation is a hyperparameter: Cherry-picked
 632 self-supervision for unsupervised anomaly detection is creating the illusion of success. *Trans. Mach. Learn.
 633 Res.*, 2023, 2023.
- 634 Sangwoong Yoon, Young-Uk Jin, Yung-Kyun Noh, and Frank Park. Energy-based models for anomaly
 635 detection: A manifold diffusion recovery approach. *Advances in Neural Information Processing Systems*,
 636 36:49445–49466, 2023.
- 637 Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In
 638 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022.
- 639 Jingyang Zhang, Jingkang Yang, Pengyun Wang, Haoqi Wang, Yueqian Lin, Haoran Zhang, Yiyou Sun,
 640 Xuefeng Du, Kaiyang Zhou, Wayne Zhang, et al. Openood v1. 5: Enhanced benchmark for out-of-
 641 distribution detection. *arXiv preprint arXiv:2306.09301*, 2023.
- 642 Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multi-
 643 variate time series forecasting. In *ICLR*, 2023.
- 644 Yue Zhao and Leman Akoglu. Toward unsupervised outlier model selection. In *International Conference on
 645 Automated Machine Learning (AutoML)*, 2024.
- 646 Yue Zhao, Ryan Rossi, and Leman Akoglu. Automatic unsupervised outlier model selection. *Advances in
 647 Neural Information Processing Systems*, 34:4489–4502, 2021.
- 648 Yue Zhao, Sean Zhang, and Leman Akoglu. Toward unsupervised outlier model selection. In *2022 IEEE
 649 International Conference on Data Mining (ICDM)*, pp. 773–782. IEEE, 2022.
- 650 Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. AnomalyCLIP: Object-agnostic
 651 prompt learning for zero-shot anomaly detection. In *The Twelfth International Conference on Learning
 652 Representations*, 2024. URL <https://openreview.net/forum?id=buC4E91xZE>.
- 653 Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Dae ki Cho, and Haifeng Chen. Deep
 654 autoencoding gaussian mixture model for unsupervised anomaly detection. In *International Conference on
 655 Learning Representations*, 2018.

656 **Appendix**

657 **Table of Contents**

658 We detail the contents in the appendix below.

- 659 • **Appendix A. Illustration of Synthetic Sata in 2-d** illustrates the inlier and outlier data synthesis
660 for pre-training with a 2-dimensional example.
- 661 • **Appendix B. Linear Transform for Scalable GMM Data Synthesis** contains the proofs for
662 efficient data synthesis in Section 3.3.
- 663 • **Appendix C. Implementation Details** includes the training and inference details of FoMo-0D.
- 664 • **Appendix D. Detailed Experiment Setup** introduces the details of pre-training and inference
665 datasets, baselines, and their hyperparameters.
- 666 • **Appendix E. Qualitative Analysis on Sample-to-Sample Attention** visualizes the attention
667 of FoMo-0D.
- 668 • **Appendix F. Ablation Analyses** studies different design choices of FoMo-0D.
- 669 • **Appendix G. Generalization Analyses** studies the generalization ability of FoMo-0D on out-of-
670 distribution synthetic datasets, ADBench, and benchmarks.
- 671 • **Appendix H. Performance Profile Plots** presents a comprehensive comparison of different
672 methods via the cumulative distribution.
- 673 • **Appendix I. Full Results** presents the detailed metric results of FoMo-0D and the baselines,
674 including AUROC, AUPRC, and F1.
- 675 • **Appendix J. Benchmark OD Datasets** shows the details (e.g., number of samples, features) of
676 each dataset in ADBench.
- 677 • **Appendix K. Differences to Prior Work on PFNs for Tabular Data** explains the difference
678 and innovation of FoMo-0D from previous works.
- 679 • **Appendix L. Discussion** provides the summary of our work and discussions on the limitations and
680 future directions of FoMo-0D.
- 681 • **Appendix M. Reproducibility Statement** details the codebase for FoMo-0D.

682 A Illustration of synthetic data in 2-d

683 We visualize our synthetic data in Figure 4, with 3 randomly created 2-d GMMs with the number of clusters
 684 ($N = 1, 2, 3$). We choose the 80th percentile as the criterion, such that inliers are samples drawn from the
 685 GMM and within the 80th percentile, and outliers are samples drawn from the inflated GMMs and outside of
 the 80th percentile.

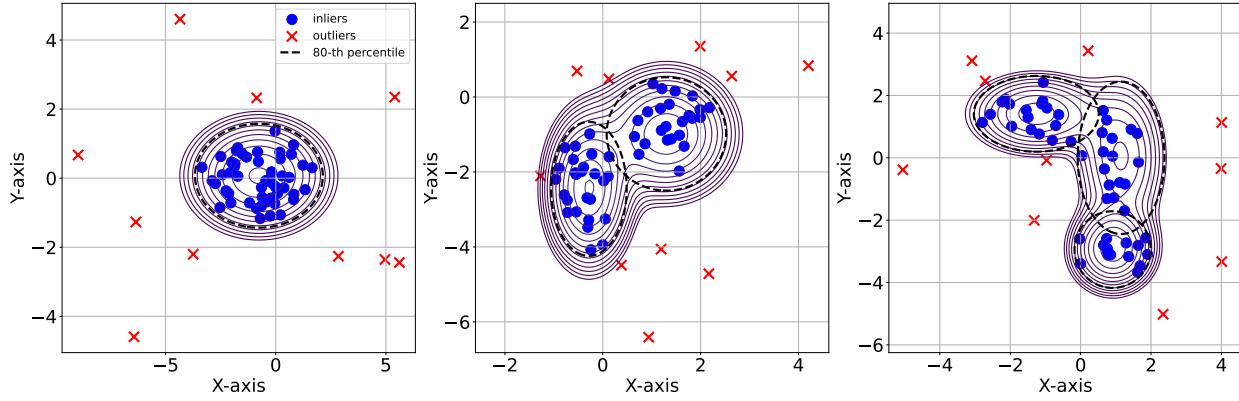


Figure 4: Illustration of synthetic data in 2D with 80th percentile as the criterion.

686

687 B Linear Transform for Scalable GMM Data Synthesis

688 B.1 Definitions

689 **Definition B.1** (Gaussian Mixture Model). We denote an m -cluster d -dimension Gaussian Mixture Model
 690 as $\mathcal{G}_m^d = \{(w_j, \mu_j, \Sigma_j)\}_{j=1}^m$, which is the weighted sum of m Gaussian distributions:

$$p(\mathbf{x}) = \sum_{j=1}^m w_j \cdot g(\mathbf{x}|\mu_j, \Sigma_j), \quad (6)$$

691 where $w_j \in \mathbb{R}^+$ is the weight for the j -th Gaussian $\mathcal{N}(\mu_j, \Sigma_j)$ with $\sum_{j=1}^m w_j = 1$, and $g(\cdot|\mu_j, \Sigma_j)$ is the
 692 density of the j -th component/cluster, with mean/center $\mu_j \in \mathbb{R}^d$ and covariance $\Sigma_j \in \mathbb{R}^{d \times d}$ being positive
 693 semi-definite, such that $\mathbf{x}^T \Sigma_j \mathbf{x} \geq 0$, for all $\mathbf{x} \in \mathbb{R}^d$.

694 **Definition B.2** (Linear Transform). We denote a linear transformation T in \mathbb{R}^d as:

$$T(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}, \quad (7)$$

695 where $\mathbf{x} \in \mathbb{R}^d$, and $\mathbf{W} \in \mathbb{R}^{d \times d}$, $\mathbf{b} \in \mathbb{R}^d$ are the parameters of T .

696 **Definition B.3** (Mahalanobis Distance). The Mahalanobis distance dist_M between a point $\mathbf{x} \in \mathbb{R}^d$ and a
 697 Gaussian distribution $\mathcal{N}(\mu, \Sigma)$ is defined as:

$$\text{dist}_M(\mathbf{x}) = \sqrt{(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)}. \quad (8)$$

698 **Definition B.4** (χ_d^2 -distribution). The Chi-squared distribution χ_d^2 with d degrees of freedom is the
 699 distribution of the sum of squares of d independent standard Normal random variables.

700 B.2 Properties

701 **Property B.5** (Lemma 5.3.2 (Casella & Berger, 2024)). If $Z \sim \mathcal{N}(0, 1)$, then $Z^2 \sim \chi_1^2$; If X_1, \dots, X_d are
 702 independent and $X_i \sim \chi_1^2$, then $\sum_{i=1}^d X_i \sim \chi_d^2$.

703 **Property B.6.** The squared Mahalanobis distance $\text{dist}_M^2(\mathbf{x}) \sim \chi_d^2$, with $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

704 *Proof:* If $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then we have $\mathbf{z} = \boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ (Gut, 2009), such that:

$$\text{dist}_M^2(\mathbf{x}) = \mathbf{z}^T \mathbf{z} = \sum_{i=1}^d z_i^2 \quad (9)$$

705 where z_i are independent standard Normal random variables. We have $\sum_{i=1}^d z_i^2 \sim \chi_d^2$ from Property B.5,
706 which completes the proof.

707 B.3 Lemmas

708 **Lemma B.7.** *Linear transform T with invertible \mathbf{W} on \mathcal{G}_m^d preserves Mahalanobis distances.*

709 *Proof:* We denote the transformed GMM as $T(\mathcal{G}_m^d) = \{(w_j, \mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}, \mathbf{W}\boldsymbol{\Sigma}_j\mathbf{W}^T)\}_{j=1}^m$, then with $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, for the transformed point $T(\mathbf{x})$ we have:

$$\text{dist}_M(T(\mathbf{x})) = \sqrt{(T(\mathbf{x}) - (\mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}))^T (\mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^T)^{-1} (T(\mathbf{x}) - (\mathbf{W}\boldsymbol{\mu}_j + \mathbf{b}))} \quad (10)$$

$$= \sqrt{(\mathbf{W}(\mathbf{x} - \boldsymbol{\mu}_j))^T (\mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^T)^{-1} (\mathbf{W}(\mathbf{x} - \boldsymbol{\mu}_j))} \quad (11)$$

$$= \sqrt{(\mathbf{x} - \boldsymbol{\mu}_j)^T \mathbf{W}^T (\mathbf{W}^T)^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{W}^{-1} \mathbf{W}(\mathbf{x} - \boldsymbol{\mu}_j)} \quad (12)$$

$$= \sqrt{(\mathbf{x} - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_j)} = \text{dist}_M(\mathbf{x}). \quad (13)$$

711

□

712 **Lemma B.8.** *Linear transform T with invertible \mathbf{W} on \mathcal{G}_m^d preserves the percentiles of the GMM.*

713 *Proof:* Let $\chi_d^2(\alpha)$ denote the α -th percentile of χ_d^2 , such that for $X \sim \chi_d^2$:

$$\text{Prob}(X \leq \chi_d^2(n)) = \frac{\alpha}{100}. \quad (14)$$

714 Based on Property B.6, we have $\text{Prob}(\text{dist}_M^2(\mathbf{x}) \leq \chi_d^2(\alpha)) = \frac{\alpha}{100}$.

715 Let $\mathbf{x} \sim \mathcal{G}_m^d$, such that $\text{dist}_M^2(\mathbf{x}) > \chi_d^2(\alpha)$ for all $\mathcal{N}_j(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, which indicates that \mathbf{x} is outside the α -th
716 percentile of \mathcal{G}_m^d . Since $\text{dist}_M(\mathbf{x})$ is preserved under T (see Lemma B.7), then we conclude that the linear
717 transform T with invertible \mathbf{W} preserves the percentiles of the GMM. □

718 C Implementation details

719 C.1 Hardware

720 We base our experiments on a NVIDIA RTX A6000 GPU with AMD EPYC 7742 64-Core Processors.

721 C.2 Training and inference

722 We train our models for 200 epochs with the Adam optimizer (Kingma & Ba, 2017) and a `learning_rate` =
723 0.001, and test with the model corresponding to the lowest training loss. The size of our $D = \{20, 100\}$
724 model is 4.87M and 4.89M parameters, respectively. We show the training process of PFNs and our model in
725 Algorithm 1.

726 **Dealing with varying dimensions and dataset size** To handle input with varying number of d features,
727 we follow Müller et al. (2022). Specifically for $d < D$, we rescale the input with $\frac{D}{d}$ and pad the features to
728 size D with 0s; and for $d > D$, we randomly sample D features out of d . In addition, FoMo-0D uses context
729 size of up to 5K at inference, where for each test sample $\mathbf{x} \in \mathcal{D}_{\text{test}}$, we randomly sample (5K–1) points as
730 $\mathcal{D}_{\text{train}}$ from datasets with $n > 5\text{K}$.

Algorithm 1: Prior-fitting of a PFN (Müller et al., 2022) and ours

Input : A prior distribution over datasets $p(\mathcal{D})$, from which samples can be drawn and the number of datasets Q to draw for one epoch, the number of training epochs E , the periodicity P , the number of unique datasets q , linear transformation T .

Output : A model q_θ that will approximate the PPD

```

1 Initialize the neural network  $q_\theta$ ;
2 Initialize the epoch-level collection  $\mathcal{C}_E = [ ]$ ;
3 for  $i \leftarrow 1$  to  $E$  do
4   if  $i \leq P$  then
5     Initialize an empty buffer  $\mathcal{B}_i = [ ]$ ;
6     Initialize the dataset-level collection  $\mathcal{C}_q = [ ]$ ;
7     for  $j \leftarrow 1$  to  $Q$  do
8       if  $j \leq q$  then
9         Step 1: sample  $D_j := \mathcal{D}_{\text{train}} \cup \{(\mathbf{x}_k, y_k)\}_{i=k}^{|\mathcal{D}_{\text{test}}|} \sim p(\mathcal{D})$ ;
10         $\mathcal{C}_q \leftarrow \mathcal{C}_q + [D_j]$ 
11      end
12    else
13       $j \leftarrow j \bmod q$ 
14       $D_j \leftarrow T(\mathcal{C}_q[j])$ 
15    end
16    Step 2: compute stochastic loss approximation  $\bar{\ell}_\theta = \sum_{k=1}^{|\mathcal{D}_{\text{test}}|} (-\log q_\theta(y_k | \mathbf{x}_k, \mathcal{D}_{\text{train}}))$ ;
17    Step 3: update parameters  $\theta$  with stochastic gradient descent on  $\nabla_\theta \bar{\ell}_\theta$ ;
18     $\mathcal{B}_i \leftarrow \mathcal{B}_i + [D_j]$ 
19  end
20   $\mathcal{C}_E \leftarrow \mathcal{C}_E + [\mathcal{B}_i]$ 
21 end
22 else
23    $i \leftarrow i \bmod P$ 
24    $\mathcal{B}_i \leftarrow \mathcal{C}_E[i]$ 
25   for  $j \leftarrow 1$  to  $Q$  do
26      $D_j \leftarrow T(\mathcal{B}_i[j])$ 
27     Perform Step 2 and Step 3
28   end
29 end
30 end

```

731 **Model architecture** We use a 4-layer Transformer with hidden dimension $\text{h_dim} = 256$, a linear embedding
 732 layer at the input ($\mathbb{R}^D \rightarrow \mathbb{R}^{\text{h_dim}}$), and a 2-layer MLP layer at the output ($\mathbb{R}^{\text{h_dim}} \rightarrow \mathbb{R}^2$) for inlier vs. outlier
 733 binary classification. For each Transformer layer, we use `num_head` = 4 for each attention module and $R = 500$
 734 for the router-based attention (Figure 2).

735 **Training loss** In Figure 5, we plot the training loss of our $D = 100$ model trained with 8K unique
 736 datasets/epoch (denoted as “8K”) versus 0.5K unique + 7.5K transformed datasets/epoch (denoted as
 737 “0.5K+T”), together with the $D = 20$ model trained with reuse periodicity $P = 1$ (denoted as “P=1”, reusing
 738 the same 8K datasets across epochs) and $P = 1$ with transformation (denoted as “P=1+T”, transforming the
 739 8K datasets across epochs). Notice that the loss with transformation is slightly higher than no transformation
 740 (i.e., $D = 100$, “0.5K+T” vs. “8K”, and $D = 20$, “P=1+T” vs. “P=1”) across all 200 epochs, which is
 741 reasonable since the transformed datasets have non-diagonal covariances that make the learning task harder
 742 and thus result in a higher training loss. The training losses of FoMo-OD with $D = 100$ are also higher than
 743 with $D = 20$ since the subspace OD tasks are harder in higher dimensions.

744 **Inference time** Figure 10 (left) showed the inference time of FoMo-OD on CPU, comparing typical attention
 745 versus the router-based attention (with $R = 500$ routers) under varying context sizes from 1K to 10K. The

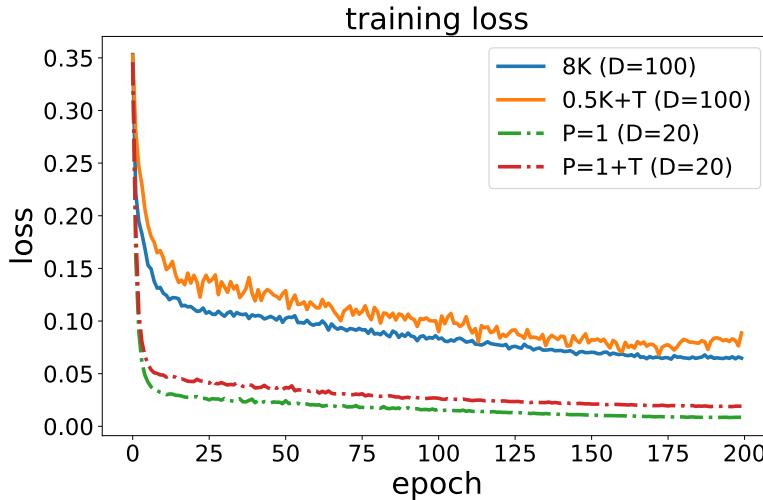


Figure 5: (best in color) Training loss of FoMo-0D ($D = 100$) with 8K unique datasets/epoch (in blue) and using 0.5K unique + 7.5K transformed datasets/epoch (in orange), and FoMo-0D ($D = 20$) with $P = 1$ (in green) and $P = 1$ with transformation (in red) over 200 epochs.

time is measured on CPU to clearly showcase the scalability trends; *quadratic* without routers and *linear* with routers.

Figure 6 shows the inference time on GPU. Notice that the time is much lower (in milliseconds), thanks to the Transformer architecture taking advantage of GPU parallelism, while the compute time for attention without routers continues to grow faster than that with routers.

In implementation, FoMo-0D (with $R = 500$ routers) uses inference context size of 5K by default, which takes about 7.7 ms per test sample on average.

D Detailed Experiment Setup

D.1 Pre-training Dataset Synthesis

During pretraining, we generate unique GMM datasets by first drawing a configuration, including dimensionality $d \in [D]$, number of components $m \in [M]$, centers $\{\boldsymbol{\mu}_j\}_{j=1}^m$ (each $\boldsymbol{\mu}_j \in [-5, 5]^d$) and covariances $\{\boldsymbol{\Sigma}_j\}_{j=1}^m$ ($\text{diag}(\boldsymbol{\Sigma}_j) \in [-5, 5]^d$). We set $M = 5$ and vary $D \in \{20, 100\}$ to study pretraining with relatively small and high dimensional datasets, respectively. We synthesize inliers and outliers as described in Section 3.1.

We then sample $S = 5,000$ points that are within the 90th percentile of the GMM. To synthesize outliers, we “inflate” a *subset* of dimensions by randomly choosing $|\mathcal{K}| \in [D]$ dimensions and multiplying the corresponding variances by $\times 5$ (following (Han et al., 2022)), i.e. $5 \times \boldsymbol{\Sigma}_{j,kk}$ ’s for $k \in \mathcal{K}$, and then draw $S = 5,000$ samples from the inflated GMM that are outside the 90th percentile of the original GMM.

To speed up data synthesis via linear transformations, we first draw 500 unique datasets using $m \in [5]$ and $d \in \{1, 2, \dots, 100\}$ (i.e. 5×100) and transform each one $15 \times$ using varying parameters (\mathbf{W}, \mathbf{b}) as described in Section 3.3.⁷ This yields 8K unique datasets (500 original and 7,500 transformed) to use at one training epoch (over 1,000 steps with batch size $B = 8$). We repeat this process at each epoch, drawing 500 new datasets and transforming them to reach 8K datasets per epoch.

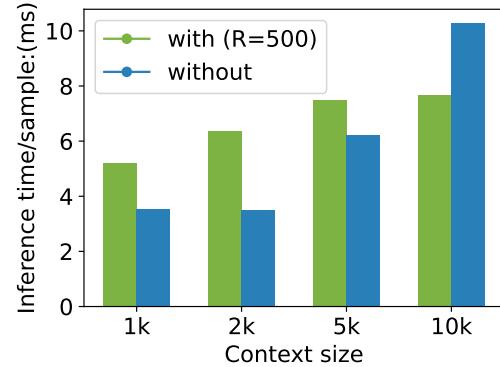


Figure 6: Inference time of FoMo-0D on *GPU* with vs. w/out router-based attention under varying context size.

⁷It is important to ensure that the eigenvalues of \mathbf{W} (i.e. variances) are not too small such that the dataset does not flatten in any direction. To this end, we draw a random orthonormal basis $\mathbf{U} \in [-1, 1]^{d \times d}$ and a diagonal $\mathbf{\Lambda}$ with eigenvalues $\lambda_{kk} \in ([-1, -0.1] \cup [0.1, 1])^d$, and obtain $\mathbf{W} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$. We also use $\mathbf{b} \in [-1, 1]^d$.

771 **D.2 Real-world Benchmark Datasets**

- 772 While pretraining is purely on synthetic datasets, we evaluate FoMo-0D on **57** real-world datasets from the
 773 ADBench benchmark (Han et al., 2022) (see Table 20). They consist of 47 popular tabular outlier detection
 774 datasets, as well as 10 newly-constructed tabular datasets created from images and natural language tasks
 775 by using pretrained models to extract embeddings. We defer to the original paper for the details on these
 776 benchmark datasets.
- 777 We compare to DTE (Livernoche et al., 2024) and baselines therein as described next, thus, following their
 778 OD setting with inlier-only $\mathcal{D}_{\text{train}}$, we split each dataset five times into train/test using five different seeds
 779 and report the mean performance and its standard deviation. In particular, each random split designates
 780 50% of the inliers as $\mathcal{D}_{\text{train}}$, while $\mathcal{D}_{\text{test}}$ contains the rest of the inliers and all the outlier samples. Note that
 781 while the baseline methods require model re-training and inference for each $\mathcal{D}_{\text{train}}/\mathcal{D}_{\text{test}}$ split, FoMo-0D uses
 782 the splits only for inference as $\mathcal{D}_{\text{train}}$ is merely passed as context.
- 783 Before passing the datasets as input to FoMo-0D, we perform a quantile transform such that the features
 784 follow a Normal distribution, to better align with the pretraining data from GMMs.

785 **D.3 Baselines**

- 786 We compare FoMo-0D against **26** baselines, from classical/shallow methods to modern/deep models. Our
 787 baselines include all the baselines imported from one of the latest papers that proposed the SOTA diffusion-
 788 based model DTE (Livernoche et al., 2024), and its three variants; DTE-C, DTE-IG, and DTE-NP. Their
 789 baselines comprise all those in ADBench (Han et al., 2022); both classical ones (k NN (Ramaswamy et al.,
 790 2000), LOF (Breunig et al., 2000), iForest (Liu et al., 2008), HBOS (Goldstein & Dengel, 2012), etc.) and
 791 deep models (DeepSVDD (Ruff et al., 2018), DAGMM (Zong et al., 2018), DROCC (Goyal et al., 2020), etc.).
 792 They also include more recent approaches based on self-supervised learning (GOAD (Bergman & Hoshen,
 793 2020), ICL (Shenkar & Wolf, 2022), SLAD (Xu et al., 2023), etc.), besides the four additional generative
 794 baselines: normalizing planar flows (Rezende & Mohamed, 2015), DDPM (Ho et al., 2020), VAE (Kingma,
 795 2013) and GANomaly (Akcay et al., 2019). We defer to the original paper for additional details. Overall, our
 796 26 baselines consist of the most recent, SOTA approaches for OD that span a diverse family (nonparametric,
 797 self-supervised, generative, etc.).

798 **D.4 Hyperparameters for Baselines**

- 799 Table 4 gives the list of HP values we used to study the HP sensitivity/performance variability of the (from
 800 top to bottom) top-4 baselines.

Table 4: Top-4 baselines (from top to bottom) and hyperparameter (HP) configurations.

Baseline	Hyperparameters
DTE-NP	$k \in \{5, 10, 20, 40, 50\}$
k NN	$k \in \{5, 10, 20, 40, 50\}$
ICL	<code>learning_rate</code> $\in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$
DTE-C	$k \in \{5, 10, 20, 40, 50\}$

801 **D.5 Ranking the 26 baselines**

- 802 Figure 24 presents the visualization of the p -values of the pairwise Wilcoxon signed rank test w.r.t. AUROC
 803 among the baseline methods used by Livernoche et al. (2024). We rank these 26 baselines based on their
 804 mean p -value (i.e., row-wise average) against the other baselines.

805 **D.6 Comparison of top-4 baseline variants with varying HP configurations**

806 Figure 25, 26, 27, 28 give the p -values, respectively comparing the variants of the top-4 baselines (DTE-NP,
 807 k NN, ICL, DTE-C) among themselves using different HP configurations, as well as the avg model with the
 808 average performance across HPs. (Specifically for ICL, `learning_rate` (`lr`) $\in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$;
 809 and for others, #nearest-neighbors $k \in \{5, 10, 20, 40, 50\}$). We find that for ICL, $\text{lr} = 10^{-3}$ or 10^{-4} are
 810 preferable while those that are too small or too large perform poorly. For others, small $k \in \{5, 10\}$ tend to
 811 outperform larger $k \in \{40, 50\}$. Note that Livernoche et al. (2024) used $k = 5$ in their paper that proposed
 812 DTE (and variants) as well as the k NN baseline for fair comparison, while the DTE^{avg} and $k\text{NN}^{\text{avg}}$ models
 813 across HP configurations perform subpar.

814 **D.7 Sampling time of d -dimensional GMM**

815 Figure 7 shows the sampling time of drawing 10,000 points
 816 from different GMMs with increasing dimensionality $d =$
 817 $\{10, 20, \dots, 200\}$. We parallelize the sampling process over
 818 10 CPUs, where each CPU draws 1000 samples.

819 We observe that the sampling time grows nonlinearly as
 820 the number of dimensions increases, which suggests that it
 821 may incur considerable computational overhead to directly
 822 draw from the data prior over hundreds of thousands of
 823 training steps, motivating the use of our proposed on-the-
 824 fly linear transformation T for scalability.

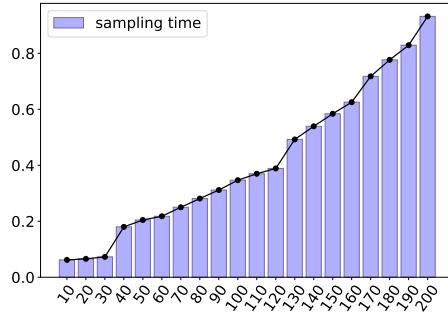


Figure 7: Sampling time (in seconds) of 10,000 points from GMMs with varying number of dimensions.

825 **E Qualitative Analysis on Sample-to-Sample Attention**

826 We sample 50 inliers as context and 100 outliers from a 2-d
 827 GMM using the 80th percentile as the labeling threshold,
 828 and visualize the top 5 inliers most attended by the 100
 829 outliers based on the average (cross) attention weights over
 830 4 heads from the last layer of FoMo-0D ($D = 100$), which
 831 accurately labeled all the 100 outliers. In Figure 8, the most
 832 frequently attended inliers are close to either the center of
 833 a Gaussian (e.g., 1st, 5th) or the criterion (e.g., 3rd, 4th),
 834 suggesting FoMo-0D tends to learn decision boundaries that
 835 reflect the prior data generation process.

836 For each outlier, we compute the sum of L2 distances to its
 837 top-5 attended inliers (`att`), the sum of L2 distances to 5
 838 randomly chosen inliers (`rdm`), and the sum of L2 distances to
 839 top-5 inliers with highest likelihood under the GMM (`prob`).
 840 We perform Wilcoxon signed rank test between `att` and `rdm`
 841 (alternative: “less”), `att` and `prob` (alternative: “greater”)
 842 over all the outliers, with a p -value of 4.4×10^{-4} and 0.99, re-
 843 spectively, suggesting the distances based on attention weights
 844 are significantly less than the random distances, and **not** significantly greater than the distances to inliers in
 845 high probability region.

846 We visualize the top-5 attended inliers for 3 outliers at different position of the 2-d GMM in Figure 9. For a
 847 specific outlier, there is a similar trend of attending to the center of a Gaussian (as shown in Figure 8), besides,
 848 inliers that reflect the criterion boundary or are close to the outlier are actively attended (e.g., 3rd, 4th in the
 849 left, 1st in the middle, 2nd, 5th in the right), suggesting FoMo-0D is incorporating both boundary and nearest
 850 neighbor information dynamically for each outlier.

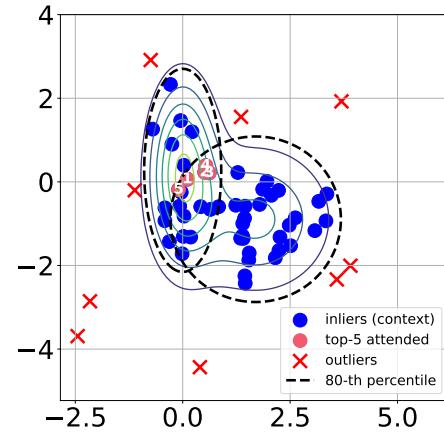


Figure 8: Top-5 attended inliers (all 50 inliers and only part of the outliers are shown for better visualization).

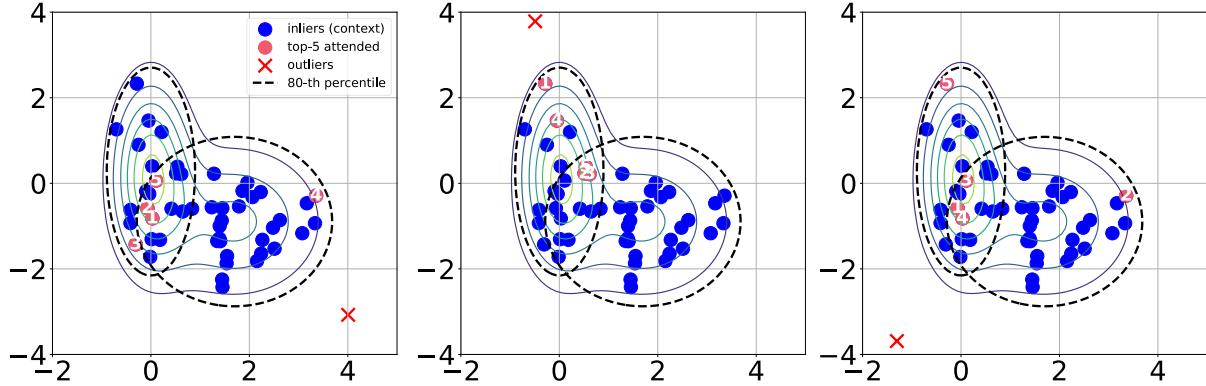


Figure 9: Top-5 attended inliers of 3 outliers at different positions of the GMM

851 F Ablation Analyses

852 In this section, we perform various ablations to study the effect of different design choices in FoMo-0D; namely,
 853 **F.1** maximum pretraining data dimensionality D , the number of routers R on **F.2** cost and **F.3** performance,
 854 **F.4** context size (both for training and inference), **F.5** number of unique datasets used for pretraining (i.e.,
 855 reuse periodicity P), data transformation T during synthesis on **F.6** performance and **F.7** speed up, **F.8**
 856 data diversity and prolonged training, **F.9** quantile transforming the benchmark datasets preceding inference,
 857 and finally, **F.10** how different model sizes affect performance.

858 Unless stated otherwise, most ablation results are performed using FoMo-0D with $D = 20$, as it is faster to
 859 pretrain under these many varying settings.

860 F.1 Effect of pretraining dimensionality D

861 *How does FoMo-0D’s generalization performance change by increasing dimensionality of the* *862 **pretraining data?***

863 We start by comparing FoMo-0D pretrained on datasets with up to $D = 20$ versus $D = 100$ dimensions. Note
 864 that learning on higher dimensional datasets is harder, as evident from the relatively larger pretraining loss
 865 as shown in Appendix Figure 5. While the statement is accurate in general, it is also partly because subspace
 866 outliers “hide” better in higher dimensions.

867 Comparing Table 1 ($D = 100$) with Table 2 ($D = 20$) w.r.t. p -values over All datasets, we find that FoMo-0D
 868 at larger scale does better, where **all** p -values are larger for $D = 100$ than $D = 20$. We find that FoMo-0D
 869 with $D = 20$ performs well on datasets with $d \leq 20$ (i.e., “on its own game”), however beyond its pretraining
 870 setting, e.g. on datasets with $d \leq 50$, $D = 100$ is superior to $D = 20$ as shown in Appendix Table 13.

871 F.2 Effect of routers on cost

872 *What is the running time and memory cost of FoMo-0D with & w/out router-based attention?*

873 Figure 10(left) shows the average inference time per test sample, comparing FoMo-0D using a router-based
 874 attention mechanism with $R = 500$ routers (in green) versus FoMo-0D using typical attention without any
 875 routers (in blue). As inference context size increases, running time for traditional attention grows quadratically
 876 while router mechanism scales linearly.⁸

877 Similarly, memory cost with routers is considerably lower when using routers, especially for larger context
 878 sizes, as shown in Figure 10(middle).

⁸Note that the inference time is reported on CPUs to show scalability. On GPUs, w/ 5K context size, see Appendix Figure 6, where typical attention takes advantage of parallelism (6.5ms), while router-based attention is slightly slower (7.7 ms w/ 500 routers) due to its **two** sequential self-attentions; see Eq.s (4) and (5).

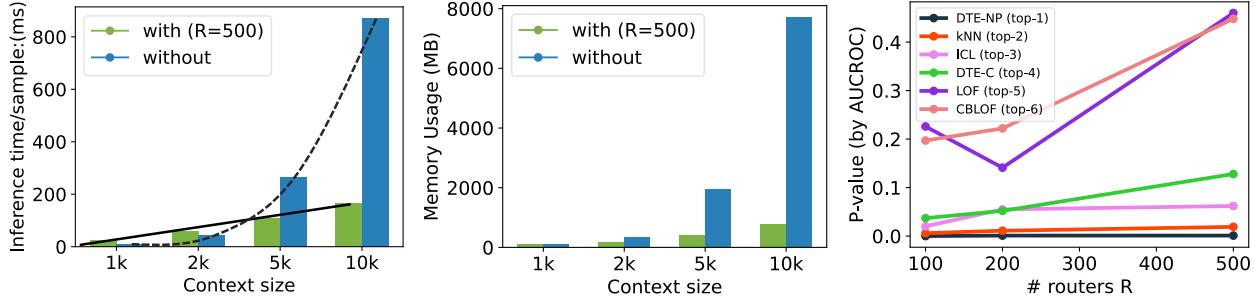


Figure 10: FoMo-0D w/ router mechanism saves time and memory while more #routers perform better, offering a cost-performance trade-off: (left) inference-time (ms) per sample and (middle) memory cost (MB) with & w/out routers by varying context size; (right) performance (based on p -value against top baselines, higher is better) vs. number of routers. (setting: $D = 20$, $P = 1$)

879 F.3 Effect of routers on performance

880 What is the impact of the number R of routers (or representatives) on performance?

881 Router-based mechanism allows to trade-off running time with expressiveness of the attention and hence
 882 performance. Figure 10(right) shows the p -values of the Wilcoxon signed rank test as the number of routers
 883 R is increased from 100 to 200 and 500, comparing FoMo-0D to each of the top-6 baselines. We notice that
 884 FoMo-0D performance tends to increase monotonically with more routers.

885 F.4 Effect of context size

886 What is the impact of context size, both during model pretraining as well as during inference?

887 To study how performance changes by context size, we train FoMo-0D with varying context size in {1K,2K,5K}
 888 and employ each pretrained model for inference with varying context size in {1K,2K,5K,10K}. Table 5 shows
 889 the results, where performance is depicted by the average rank of FoMo-0D (the lower, the better).

Table 5: Average rank (based on comparison to 30 baselines w.r.t. AUROC) of FoMo-0D across datasets under different context sizes for training and inference. Smaller ranks imply better performance. (setting: $D = 20$, $R = 500$, $P = 1$)

	Infer:1K	Infer:2K	Infer:5K	Infer:10K
Train:1K	13.816	14.623	15.193	15.439
Train:2K	13.079	13.219	13.439	13.561
Train:5K	13.088	13.211	13.307	13.430

890 We find that training with a larger context improves performance at any inference context size. On the other
 891 hand, perhaps counter-intuitively, FoMo-0D with smaller inference context size does better. We conjecture
 892 that is because the #routers-to-context size ratio increases with a larger context size at inference, limiting
 893 the expressive power of the “bottleneck” attention mechanism. The pairwise statistical tests among the
 894 $3 \times 4 = 12$ models support these observations, as shown in Figure 11. Interestingly, when the training context
 895 size is large enough at 5K, inference with 10K samples generalizes beyond training with no statistical evidence
 896 for performance difference (at 0.05) from other inference context sizes.

897 F.5 Effect of number of unique datasets

898 How do FoMo-0D performances compare when pretrained on unique vs. reused datasets, via 899 varying periodicity P ?

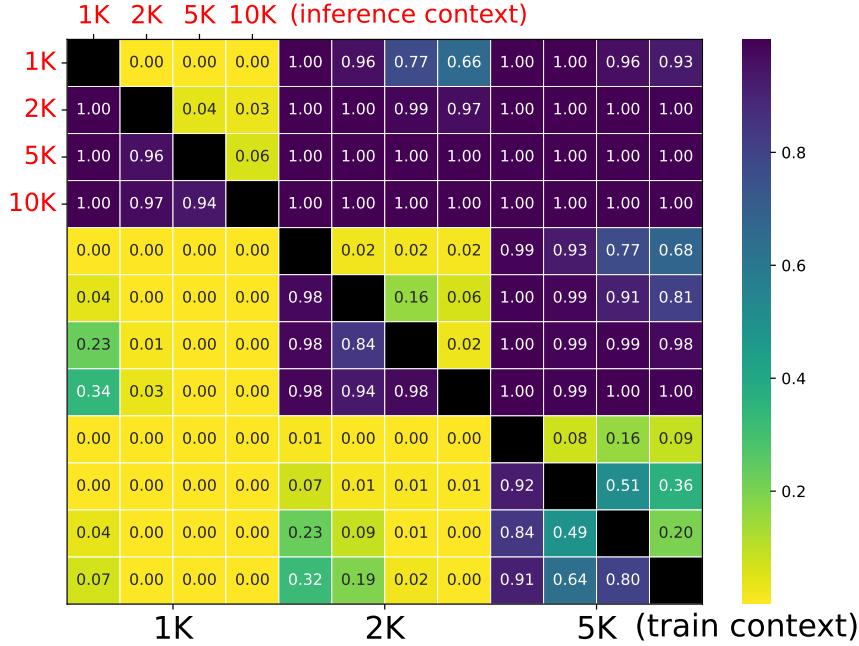


Figure 11: p -values of the pairwise Wilcoxon signed rank test between models (larger p implies col-method is better than row-method) w/ different context sizes for **training** (1K/2K/5K, 1st/2nd/3rd four grids, in **black**) and **inference** (1K/2K/5K/10K, every 1st/2nd/3rd/4th grid, in **red**): Larger training context improves overall performance, while smaller inference context is preferable.

Table 6: Ablation results on dataset reuse across epochs with varying $P \in \{1, 50, 100\}$ show stable p -values against the top-5 baselines, where there is no statistical evidence to suggest performance difference between FoMo-0D with $D = 20$ and the top 3rd baseline at 0.05 w.r.t. pairwise Wilcoxon signed rank test comparisons, while it continues to significantly outperform the top 5th baseline (LOF) when $d \leq 50$. (setting: $D=20$, $R=500$, context size=5K, w/out transformation T)

	$P = 1$ (#unique datasets: 8K)					$P = 50$ (#unique datasets: $8 \times 50 = 400$ K)					$P = 100$ (#unique datasets: $8 \times 100 = 800$ K)				
top-5	DTE-NP	kNN	ICL	DTE-C	LOF	DTE-NP	kNN	ICL	DTE-C	LOF	DTE-NP	kNN	ICL	DTE-C	LOF
All	<u>0.001</u>	<u>0.019</u>	0.062	0.128	0.460	<u>0.001</u>	<u>0.019</u>	0.089	0.159	0.394	<u>0.001</u>	<u>0.015</u>	0.072	0.121	0.290
$d \leq 20$	0.583	0.755	0.943	0.736	0.998	0.572	0.789	0.968	0.616	0.993	0.439	0.678	0.953	0.550	0.972
$d \leq 50$	0.415	0.750	0.869	0.962	0.999	0.347	0.794	0.893	0.946	0.997	0.293	0.697	0.890	0.924	0.994

900 Next we study the effect of dataset *reuse at epoch level* (w/out transformation) on performance as presented
901 in Section 3.3. We vary reuse periodicity P in $\{1, 50, 100\}$, and accordingly, increase the number of unique
902 datasets used for pretraining across epochs. As shown in Table 6, FoMo-0D (w/ $D = 20$) performs similarly
903 with varying dataset reuse. In fact, it is competitive even with $P = 1$, remaining no different from the 3rd
904 best baseline (ICL) across All (57) datasets, while significantly outperforming the top 5th (LOF) across (24)
905 datasets with $d \leq 20$ as well as (38) with $d \leq 50$.

906 F.6 Effect of transformation T for synthesis

907 **How do FoMo-0D performances compare when pretrained on datasets with vs. w/out linear
908 transformation?**

909 Setting $P = 1$, we next study the impact of linear transformation T . Table 7 presents the results, where we
910 compare reuse of the *same* 8K unique datasets across epochs (w/out T), versus *transforming* these datasets
911 with T at every epoch with different parameters (w/ T). FoMo-0D performance remains stable; no statistical
912 evidence for performance difference from the top 3rd model on All datasets, while significantly outperforming
913 the top 5th across those with $d \leq 20$ and $d \leq 50$. This suggests that T can be employed without sacrificing
914 performance to save time during pretraining.

Table 7: Ablation results on performance w/ & w/out linear transformation T show stable p -values against the top-5 baselines, with no statistical evidence for performance difference between FoMo-0D with $D = 20$ and the top 3rd baseline at 0.05 w.r.t. pairwise Wilcoxon signed rank test comparisons. (setting: $D = 20$, $R = 500$, context size=5K, $P = 1$)

top-5	w/out transformation T					w/ transformation T				
	DTE-NP	kNN	ICL	DTE-C	LOF	DTE-NP	kNN	ICL	DTE-C	LOF
All	0.001	0.019	0.062	0.128	0.460	0.002	0.015	0.226	0.210	0.280
$d \leq 20$	0.583	0.755	0.943	0.736	0.998	0.648	0.708	0.988	0.718	0.955
$d \leq 50$	0.415	0.750	0.869	0.962	0.999	0.264	0.382	0.971	0.900	0.963

915 F.7 Speed up by T

916 *What is the time saving on data synthesis with linear transfor-* 917 *mation?*

918 Figure 12 shows the distribution of pretraining running-time per epoch
919 with and w/out data transformation. Specifically, we compare (left) gen-
920 erating 8K unique datasets/epoch on-the-fly and (right) first generating
921 500 unique datasets on-the-fly and then transforming each one 15 times
922 using T with different parameters to reach 8K datasets at each epoch.

923 Notice that pretraining with T takes about 450 sec./epoch on average,
924 while without T it requires 1200 sec./epoch to generate 8K unique datasets
925 and gradient descent across 1000 steps. Different from other ablation
926 results, which are based on the $D = 20$ model, here we report the running
927 times for our $D = 100$ model. Overall, our final FoMo-0D took ≈ 25 hours
928 for pre-training (450 sec. \times 200 epochs). Importantly, this is a one-time
929 cost that amortizes across many downstream tasks with as low as **7.7 ms**
930 **inference time** per test sample (see Table 3 and Appendix Figure 6).

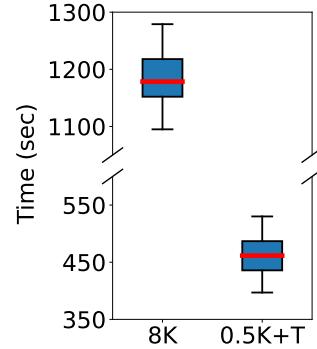


Figure 12: Runtime/epoch dist.n over 100 epochs for FoMo-0D ($D=100$) with (left) $P=100$, i.e. 8K unique datasets/epoch vs. (right) 0.5K unique+7.5K transformed datasets/epoch.

931 F.8 Effect of data diversity and prolonged training

932 *How does FoMo-0D’s performance change by increasing pretraining data diversity and number* 933 *of training epochs?*

934 Originally we have trained FoMo-0D w/ $D = 100$ using 0.5K unique + 7.5K transformed datasets over 200
935 epochs. As mentioned earlier, learning in higher dimensions tends to incur a larger loss in general but also
936 specifically here, as subspace outliers are harder to detect in high dimensions.

937 Toward reducing the loss further, we resume the pretraining for another 100 epochs. Further, to simplify the
938 tasks and thereby increase data diversity, we also decrease the inlier/outlier labeling percentile threshold from
939 90% to 80% during on-the-fly data generation in the last 100 epochs. In Figure 13, we present the training
940 loss of FoMo-0D ($D = 100$) trained with 0.5K unique + 7.5K transformed datasets/epoch over 200 epochs
941 (90th percentile as labeling threshold) and then 100 additional epochs (80th percentile as the threshold) to
942 show how data diversity and amount affect model performance.

943 Figure 14 compares FoMo-0D’s performance (w/ $D = 100$) to top-5 baselines w.r.t. p -values of the paired
944 Wilcoxon signed rank test on datasets with $d \leq 100$, after the first 200 epochs versus after 300 epochs. The
945 increase in all the p -values showcases the benefit of additional training.

946 F.9 Effect of applying quantile transform on benchmark datasets

947 *What is the impact of quantile data transform preceding inference on performance?*

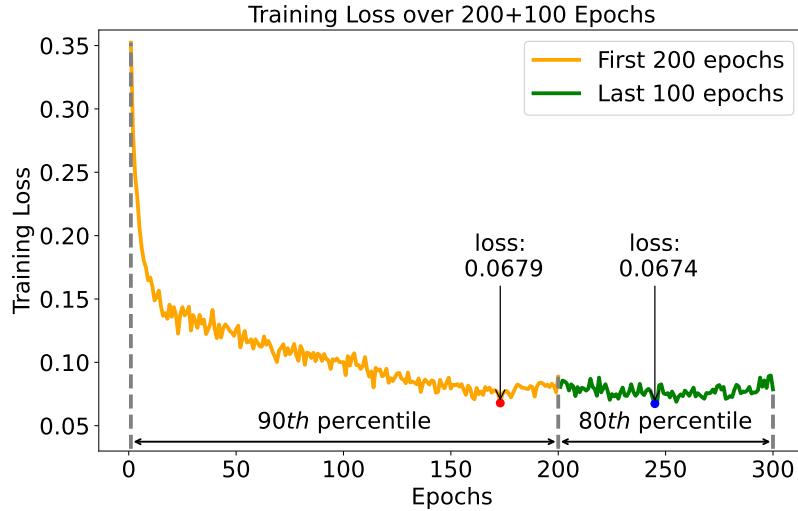


Figure 13: (best in color) Training loss of FoMo-0D ($D = 100$) with 0.5K unique + 7.5K transformed datasets/epoch for 200 epochs (in orange), followed with additional 100 epochs of training (in green). For the first 200 epochs we train with 90th percentile as the inlier/outlier threshold, which we reduce to 80th in the subsequent 100 epochs.

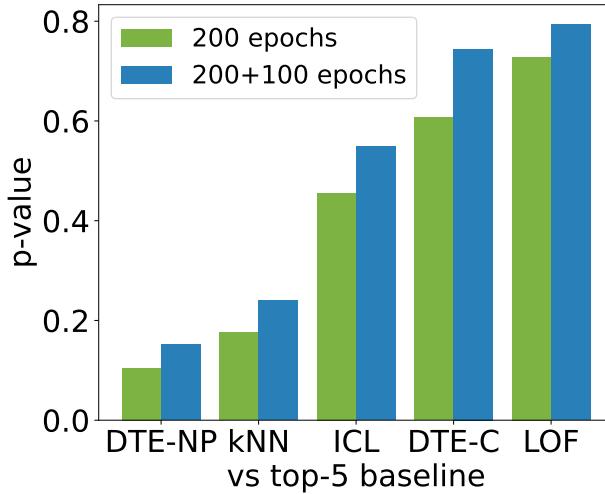


Figure 14: p -values increase with additional 100 epochs of pretraining, i.e. FoMo-0D w/ $D = 100$ performs better against top-5 baselines on datasets w/ $d \leq 100$.

948 We pretrain FoMo-0D on synthetic datasets from a simple data prior based on GMMs. The real-world
949 benchmark datasets, on the other hand, may exhibit features with distributions different from Gaussians.
950 To close the gap, we apply a quantile transform (denoted QT) on the benchmark datasets prior to feeding
951 them to FoMo-0D for inference, which transforms the features to exhibit a more Gaussian-like probability
952 distribution.

953 Figure 15 compares the performance of three FoMo-0D w/ $D = 100$ variants with and w/out QT against the
954 top-5 baselines w.r.t. the p -values of the paired Wilcoxon signed rank test. FoMo-0D tends to perform better
955 as suggested by larger p -values when QT is applied.

956

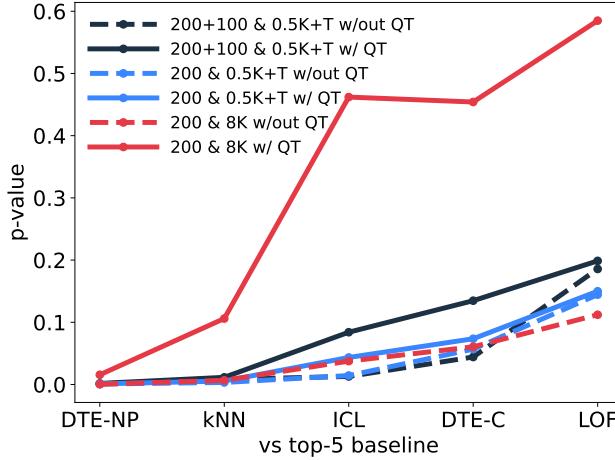


Figure 15: p -values increase, i.e. FoMo-0D performance improves, against top-5 baselines with quantile transform (QT) preceding inference, for 3 different settings of FoMo-0D w/ $D = 100$.

Table 8: p -values of the one-sided Wilcoxon signed rank test, comparing FoMo-0D (with $D = 100$) to **top 10 baselines** with default hyperparameters (HPs), and **top 4^{avg}** baselines⁶ with **avg**. performance over varying HPs (denoted w/ ^{avg}) over all (57) datasets. FoMo-0D improves (i.e. higher p -values) as number of transformer layers L increases. At $L = 1$, in-context learning ability appears to be quite limited. (setting: $D = 100$, $R = 500$, train/inference context size=5K, w/ quantile transform, $\alpha = 0.05$)

FoMo-0D	#parameters	DTE-NP	kNN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP ^{avg}	kNN ^{avg}	ICL ^{avg}	DTE-C ^{avg}
$L = 1$	1.34M	1.66×10^{-9}	4.30×10^{-9}	1.25×10^{-6}	1.34×10^{-7}	5.08×10^{-6}	5.44×10^{-8}	1.31×10^{-4}	1.07×10^{-6}	1.30×10^{-6}	1.46×10^{-7}	2.68×10^{-9}	1.53×10^{-8}	2.13×10^{-7}	0.395
$L = 2$	2.52M	0.006	0.036	0.157	0.259	0.442	0.557	0.703	0.431	0.759	0.805	0.021	0.134	0.333	1.000
$L = 3$	3.70M	0.016	0.098	0.372	0.579	0.572	0.871	0.808	0.652	0.921	0.961	0.085	0.335	0.652	1.000
$L = 4$	4.89M	0.016	0.106	0.462	0.454	0.585	0.750	0.823	0.759	0.901	0.895	0.112	0.315	0.670	1.000

957 F.10 Effect of Model Size

To understand how the performance of FoMo-0D scales with model sizes, we vary the number of transformer layers $L = 1, 2, 3$, where the default is $L = 4$ (see in **Model architecture**). We present the p -values of different model sizes in Table 8, where a p -value > 0.05 means no statistical evidence to suggest performance difference between FoMo-0D and the compared baseline. We can observe that FoMo-0D is not comparable to the top-10 baselines with one layer, and shows improved performance (i.e., p -value increases and becomes larger than 0.05) as the number of layers increases, which suggests that scaling up the model size might help improve FoMo-0D’s performance. On the other hand, we can see that the detection performance didn’t increase a lot from $L = 3$ to 4, which suggests there exist diminishing returns in scaling the model size, and to further improve performance, one has to consider other methods (e.g., increase prior complexity, more pre-training epochs) besides scaling up only the model size.

968 G Generalization Analyses

969 G.1 Generalization to Out-of-Distribution Synthetic Datasets

We conduct analyses to understand FoMo’s ability to generalize on out-of-distribution synthetic GMM datasets. Besides the in-distribution setting for pre-training (i.e., $\mu \in [-5, 5], \Sigma \in (0, 5], m \leq 5, d \leq 100$), we consider the following out-of-distribution settings: **(a)** mean and covariance significantly out of range, with $\mu \in [-50, -5] \cup [5, 50], \Sigma \in [5, 50]$, denoted as “ $|\mu|, |\Sigma| \in [5, 50]$ ”; **(b)** number of clusters significantly out of range, denoted as “ $m \in [5, 50]$ ”; **(c)** number of dimensions significantly out of range, denoted as “ $d \in [100, 500]$ ”; **(d)** binary outliers with values either 0 or 1 in one dimension from the sub-dimensions, denoted as “binary”; **(e)** “all”, which combines all the variants above. For each setting, we generate 1000 datasets with random seeds from 0 to 999, where on each dataset, we simulate 1000 test points with an

978 outlier rate of 5% and evaluate FoMo-0D with a context length of 5000. We present the results with averaged
 979 performance over 1000 datasets for each setting in Table 9.

Table 9: Average metric score \pm standard dev. over 1000 seeds for different out-of-distribution (OOD) synthetic GMMs. FoMo-0D remains robust against OOD test datasets as in (a)–(d), maintaining similar performance to in-distribution performance (top). Performance is affected more when datasets are OOD w.r.t. multiple factors combined as in (e).

Dataset	AUROC	AUCPR	F1
ID: in-distribution	98.55 ± 2.73	91.17 ± 13.07	86.74 ± 15.43
(a) OOD w.r.t. $ \mu , \Sigma \in [5, 50]$	94.79 ± 7.53	80.62 ± 21.19	76.32 ± 19.85
(b) OOD w.r.t. $m \in [5, 50]$	97.69 ± 3.59	86.72 ± 15.57	81.20 ± 16.23
(c) OOD w.r.t. $d \in [100, 500]$	96.22 ± 9.01	86.37 ± 23.27	83.23 ± 22.08
(d) OOD w.r.t. binary variable	100.00 ± 0.00	100.00 ± 0.06	99.97 ± 0.34
(e) OOD w.r.t. all combined	85.44 ± 16.96	64.17 ± 35.07	63.99 ± 33.53

980 We can observe different extents of performance degradation when applying out-of-distribution variations.
 981 Compared to other single variations, FoMo-0D seems to suffer more from inflating the mean and covariances,
 982 as due to the significant deviation in the parameters of the GMMs, inliers generated under such a setting are
 983 seemingly “outliers” w/o any reference points. Surprisingly, although FoMo-0D is only trained on continuous
 984 data, it can almost perfectly classify binary outliers hidden in one of the sub-dimensions, suggesting FoMo-0D
 985 could potentially generalize to discrete data at test time.
 986 However, with all variations added, FoMo-0D becomes less capable compared to one single out-of-distribution
 987 variation, although there might exist some signals (e.g., binary labels) in favor of its decision-making process,
 988 for which training a powerful model with more comprehensive priors could possibly alleviate the issue.

989 G.2 Generalization to Out-of-GMM-Distribution Real-World Datasets

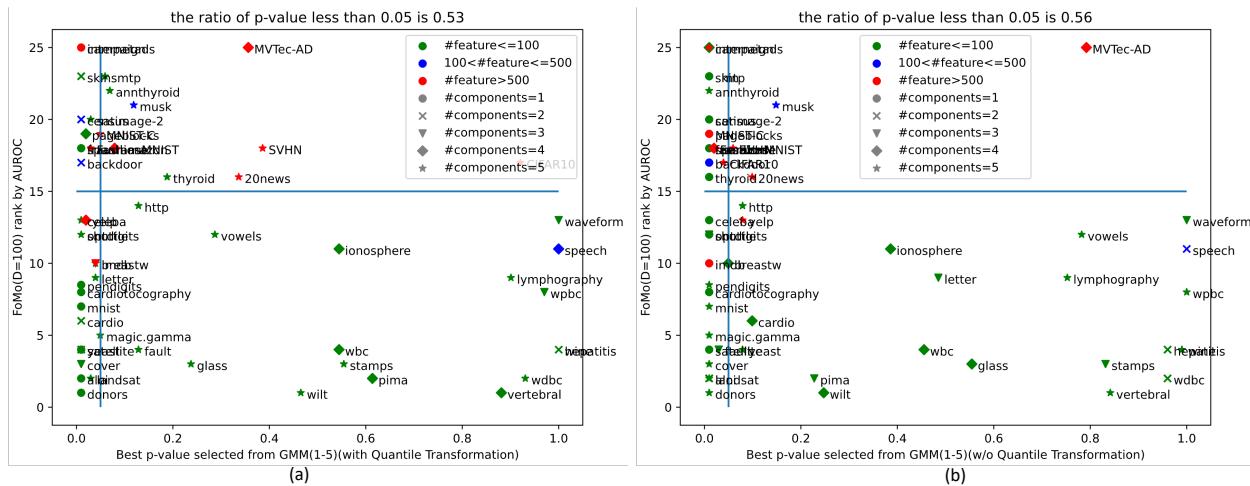


Figure 16: Goodness of fit test results for 57 datasets in AD-Bench. The x-axis is the p -value of a dataset, where a small p -value indicates GMMs are not a good fit for the dataset, and the y-axis is the rank of FoMo-0D ($D=100$) on that dataset (the smaller the better). We plot p -value = 0.05 (vertical) and rank = 15 (horizontal) as references. The left figure (a) is with quantile transformation while figure (b) is without quantile transformation. We use different colors to represent datasets at different dimensions, and use different markers to represent different numbers of clusters.

990 To understand how FoMo-0D generalizes from the pre-training GMM data priors to complex real-world
 991 datasets when performing zero-shot OD, we conduct the goodness of fit test Huber-Carol et al. (2012) for

992 datasets in ADBench. We fit GMMs to each real-world dataset D_{real} with up to 5 components (as with our
993 pre-training datasets), then sample D_{syn} from the best-parameter-fitted GMM and perform a two-sample
994 test⁹ on D_{real} and D_{syn} , with the null hypothesis that they come from the same distribution. A smaller
995 p -value (≤ 0.05) of such a test provides evidence toward rejecting the null, which suggests GMM is not a
996 good fit for the dataset (i.e., the pre-training data distribution is different from the test data).

997 We present the results in Figure 16, depicting the p -value (of the goodness-of-GMM-fit test) vs. FoMo-0D
998 's performance rank (the lower the better) among 30 baselines. We report the result both with quantile
999 transformation in Figure 16(a) and without quantile transformation in 16(b). Since the two figures are
1000 highly similar, our next analysis will primarily focus on the figure with quantile transformation to align with
1001 our model's implementation. We plot the vertical and horizontal lines as p -value = 0.05 and rank = 15. For
1002 p -value ≥ 0.05 and rank < 15 , we observe that performance is good on datasets with relatively large p -value
1003 where we cannot reject the null (i.e. GMM is a relatively good fit). This is where arguably FoMo-0D recalls
1004 its data prior distribution and generalizes to datasets similar to those seen during pretraining. We also see,
1005 for p -value < 0.05 and rank ≥ 15 , datasets with relatively poor performance where we can reject the null (i.e.
1006 GMM is not a good fit). These can be attributed to falling short in generalization to OOD datasets.

1007 On the other hand, we observe many datasets concentrate on p -value < 0.05 and rank < 15 , where p -value is
1008 small (GMM not a good fit) yet the performance is competitive — those are the datasets on which FoMo-0D
1009 is likely to have achieved out-of-distribution generalization. It remains an open (theoretical) question to
1010 understand what (algorithm, if any) FoMo-0D might have learned that generalizes to out-of-distribution
1011 datasets. It is also an open (empirical) quest to explore whether a more complex data prior, beyond GMMs,
1012 could further push the performance up and by how much.

1013 G.3 Generalization to Out of Distribution (OOD) Detection Tasks

1014 We further evaluate FoMo-0D on more complex datasets (e.g., ImageNet-level). Specifically, we employ
1015 OpenOOD Zhang et al. (2023), an out-of-distribution (OOD) detection benchmark, where models are trained
1016 on labeled in-distribution datasets, with K known classes, and then evaluated on out-of-distribution datasets,
1017 aiming to detect $K + 1, K + 2, \dots$ novel classes. Although OOD detection is inherently different from OD, we
1018 can construct an OD dataset from OOD datasets, treating all K class samples as inliers and the $K + 1, K + 2, \dots$
1019 OOD samples as outliers. For the in-distribution datasets, we choose ImageNet1K, which contains 1000
1020 categories of images, and ImageNet200, a subset of ImageNet1K containing 200 categories. We further choose
1021 SSB-hard, NINCO, iNaturalist, Textures, and OpenImage-O as the out-of-distribution datasets, which gives
1022 us a total number of $2 \times 5 = 10$ datasets that are ImageNet-level complex.

1023 Following Han et al. (2022), we create 10 new OD datasets from OpenOOD containing 10,000 samples with
1024 5% outliers, and use the embedding from the last average pooling layer of ResNet18 He et al. (2016) as the
1025 feature (512) for each sample. Comparing FoMo-0D with the top-4 (on our original testbed) baselines in the
1026 order of: DTE-NP, kNN, ICL, DTE-C, we follow Livernoche et al. (2024) and report mean (standard dev.)
1027 over 5 runs (seed=0/1/2/3/4) on each dataset. We present the results with in-distribution datasets being
1028 ImageNet200 and ImageNet1K in Table 10 and 11, respectively.

Table 10: Average AUROC score \pm standard dev. over five seeds for in-distribution dataset being **ImageNet200**. We use blue and green respectively to mark the top-1 and the top-2 method.

dataset	DTE-NP	kNN	ICL	DTE-C	FoMo-0D
ssb-hard	58.03 ± 0.00	58.14 ± 0.00	60.52 ± 0.25	60.74 ± 1.88	58.34 ± 1.55
ninco	53.28 ± 0.00	54.14 ± 0.00	59.56 ± 0.63	58.83 ± 1.54	55.16 ± 2.19
inaturalist	29.38 ± 0.00	29.51 ± 0.00	35.96 ± 1.10	41.77 ± 2.84	38.85 ± 3.29
textures	59.28 ± 0.00	59.91 ± 0.00	66.40 ± 0.69	70.33 ± 3.18	59.89 ± 2.07
openimageo	52.82 ± 0.00	53.79 ± 0.00	55.20 ± 0.69	59.09 ± 1.50	54.77 ± 1.19
average	50.56	51.10	55.53	58.15	53.40

⁹We use e-test from <https://www.rdocumentation.org/packages/energy/versions/1.7-11/topics/eqdist.etest>

Table 11: Average AUROC score \pm standard dev. over five seeds for in-distribution dataset being **ImageNet1K**. We use blue and green respectively to mark the top-1 and the top-2 method.

dataset	DTE-NP	kNN	ICL	DTE-C	FoMo-0D
ssb-hard	55.63 \pm 0.00	55.94 \pm 0.00	58.79 \pm 1.20	59.17 \pm 1.82	56.73 \pm 2.65
ninco	48.23 \pm 0.00	49.10 \pm 0.00	55.25 \pm 0.87	57.60 \pm 3.93	52.70 \pm 2.70
inaturalist	30.24 \pm 0.00	30.28 \pm 0.00	35.03 \pm 1.42	41.96 \pm 3.13	38.94 \pm 4.59
textures	54.38 \pm 0.00	55.43 \pm 0.00	61.30 \pm 0.95	63.10 \pm 3.72	55.18 \pm 2.92
openimagegeo	54.31 \pm 0.00	54.91 \pm 0.00	54.02 \pm 0.43	58.71 \pm 2.08	56.95 \pm 3.89
average	48.56	49.13	52.88	56.11	52.10

We further report the p -value of the Wilcoxon signed rank test between the baselines and FoMo-0D on the 10 datasets from OpenOOD, as well as on the expanded benchmark combining those 10 with our original ADBench (10+57) in Table 12. In terms of metric values, FoMo-0D performs 2nd or 3rd best across OOD datasets. p -values show that it significantly outperforms DTE-NP and kNN ($p > 0.95$, such that the p -value < 0.05 for rejecting the null hypothesis and accepting the alternative hypothesis that the “baseline-minus-FoMo-0D” gap is smaller than zero) and is no different from ICL (2nd best after DTE-C). These results demonstrate that FoMo-0D generalizes beyond OD datasets and maintains strong zero-shot OD performance on complex, ImageNet-level OOD benchmarks.

Table 12: p -value of the Wilcoxon signed rank test (alternative: “greater”) between baselines and FoMo-0D on OpenOOD and combined benchmark on AUROC. A small p -value (≤ 0.05) means that there is statistical evidence for the alternative hypothesis such that baselines achieve higher metric performance than FoMo-0D.

method	DTE-NP	kNN	ICL	DTE-C
OpenOOD	1	0.9951	0.1875	0.0009
OpenOOD+ADBench	0.1271	0.3308	0.3153	0.1265

Interestingly, we observe that ICL and DTE-C outperform DTE-NP and kNN on the OpenOOD datasets, whereas on ADBench, DTE-NP and kNN are the top-2 methods outperforming ICL and DTE-C. We hypothesize this is because it is harder for non-parametric methods like DTE-NP and kNN to estimate meaningful decision boundaries in high dimensions (e.g., 512). In contrast, the performance of FoMo-0D is consistently competitive, where the p -values on the combined testbed (OpenOOD+ADBench) show that FoMo-0D is as competitive as all the top baselines across 67 diverse datasets, while maintaining zero-shot detection ability.

H Performance Profile Plots

To enable a comprehensive comparison of different methods, we plot rank (w.r.t. AUROC, lower is better) distribution in Figure 17, and adopt τ performance profile plots as described in Dolan & Moré (2002). These plots display the cumulative distribution of the τ metric—which quantifies suboptimality relative to the best-performing method. By computing sorted τ values along with their cumulative probabilities, we then use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

Figure 18, Figure 19, and Figure 20 illustrate performance profile plots of FoMo-0D and other baselines across all datasets. The results show that FoMo-0D (D=100) ranks at **top-5 (w.r.t. AUROC)**, **top-3 (w.r.t. AUPR)** and **top-1 (w.r.t. F1)**, respectively, outperforming many baselines.

Moreover, the performance of FoMo-0D (D=100) is even better (i.e., ranked within top-2) when tested on datasets with dimensions less than 100. As shown in Figure 21, Figure 22, and Figure 23, the area under the curve of FoMo-0D (D=100) ranks at **top-1 (w.r.t. AUROC)**, **top-2 (w.r.t. AUPR)** and **top-2 (w.r.t. F1)**, respectively, under this setting.

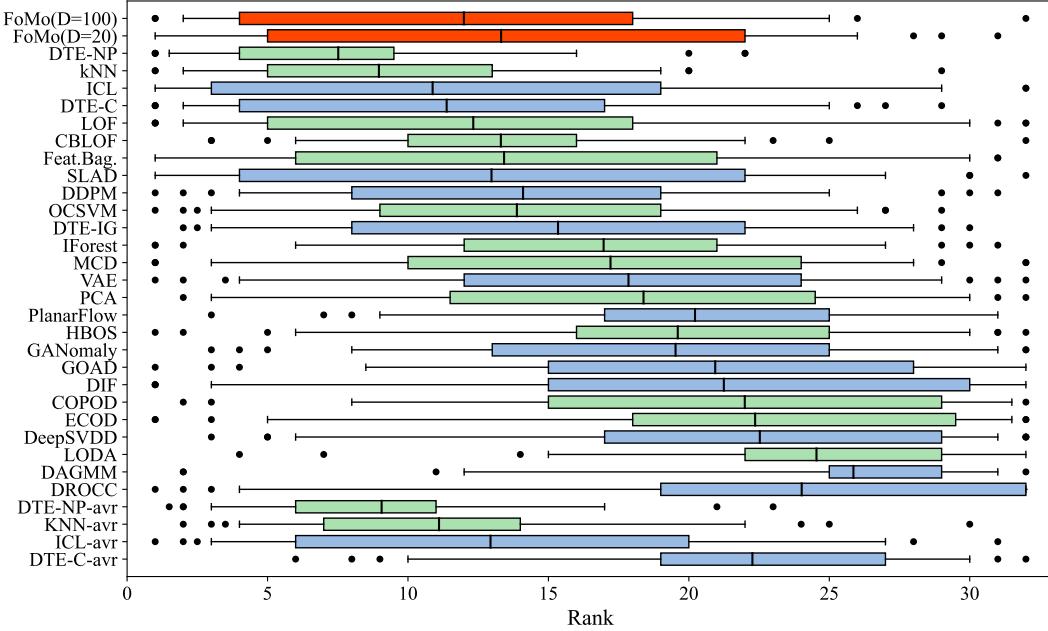


Figure 17: (best in color) Rank (w.r.t. AUROC, lower is better) distribution across all **57** real-world datasets shown via boxplots for (from top to bottom) FoMo-0D in red, all **26** baselines ordered by mean p -value⁶ (shallow and deep baselines in green and blue), and **top 4** baselines' avg variants. The vertical line depicts the mean, the box shows the 25-75%, bars range 5-95%, and circles show the datasets at the tails.

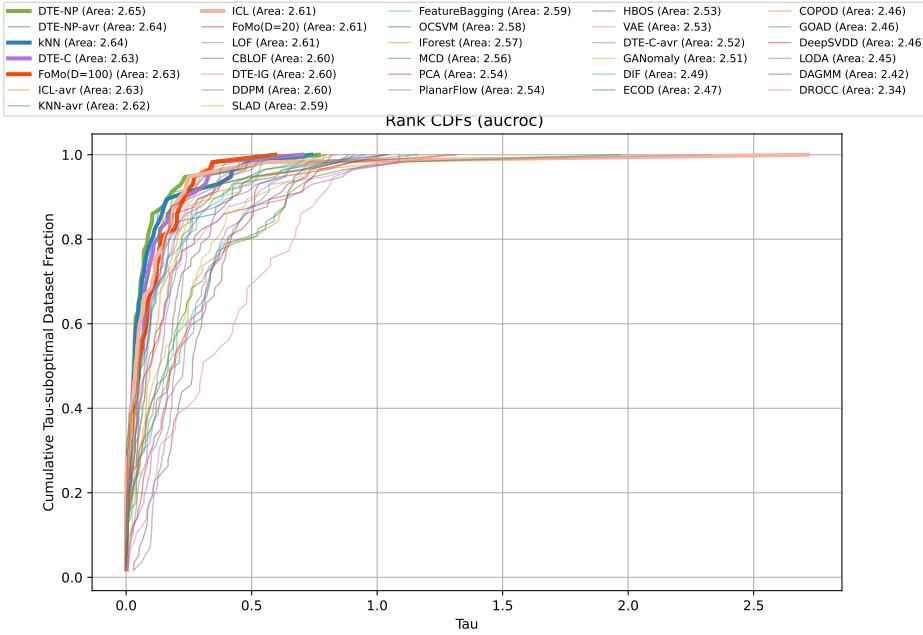


Figure 18: FoMo-0D ranks in **top-5** based on the performance profile plots of all detectors w.r.t. **AUROC** across **all datasets**. In the plot, x-axis represents the τ values—performance ratios that compare each method's metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method's performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

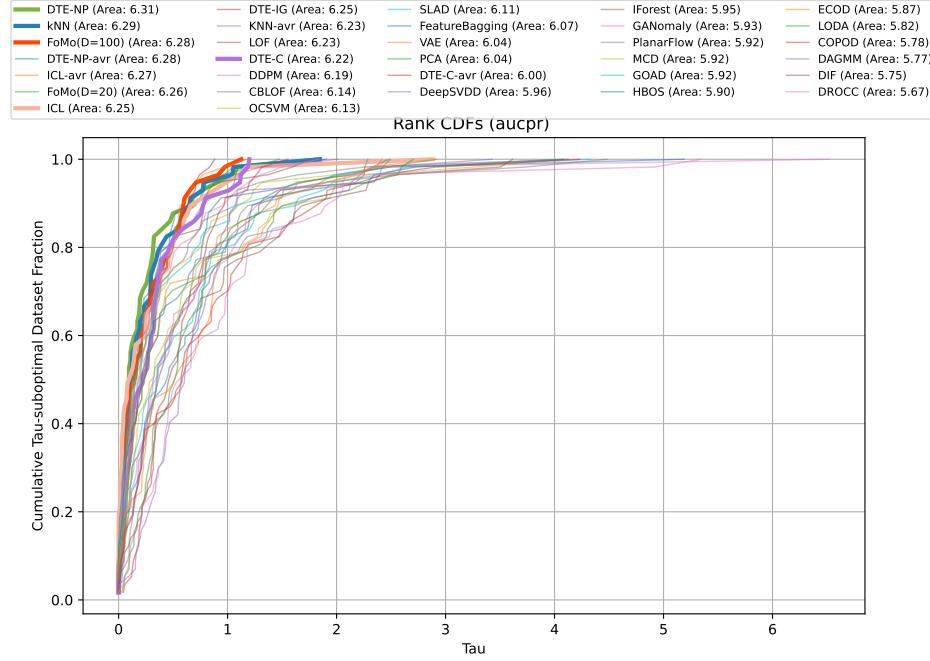


Figure 19: **FoMo-0D** ranks at **top-3** based on the performance profile plots of all detectors w.r.t. **AUPR** across **all datasets**. In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

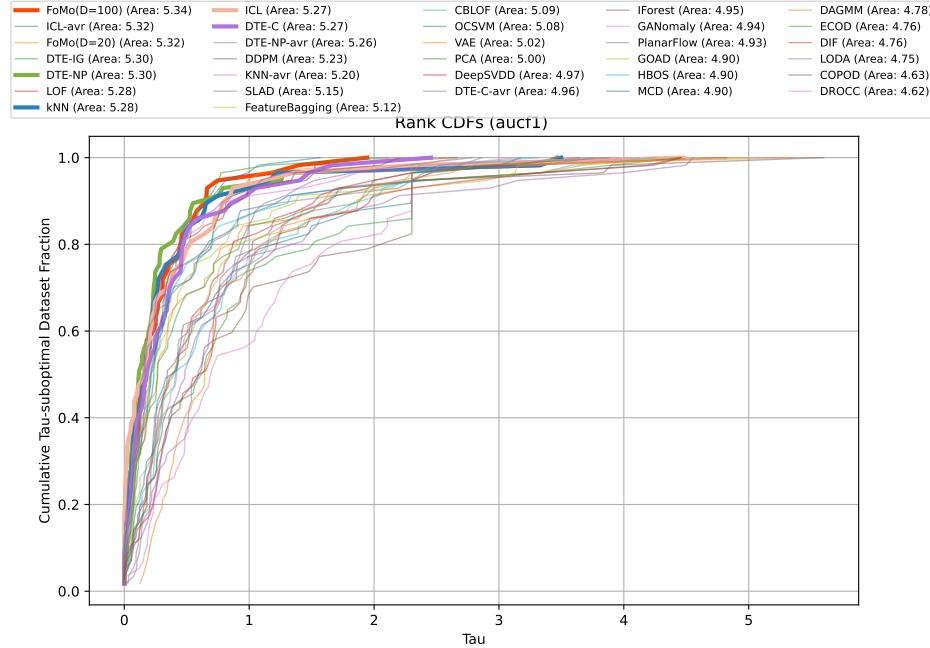


Figure 20: **FoMo-0D** ranks at **top-1** based on the performance profile plots of all detectors w.r.t. F1 across all datasets. In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

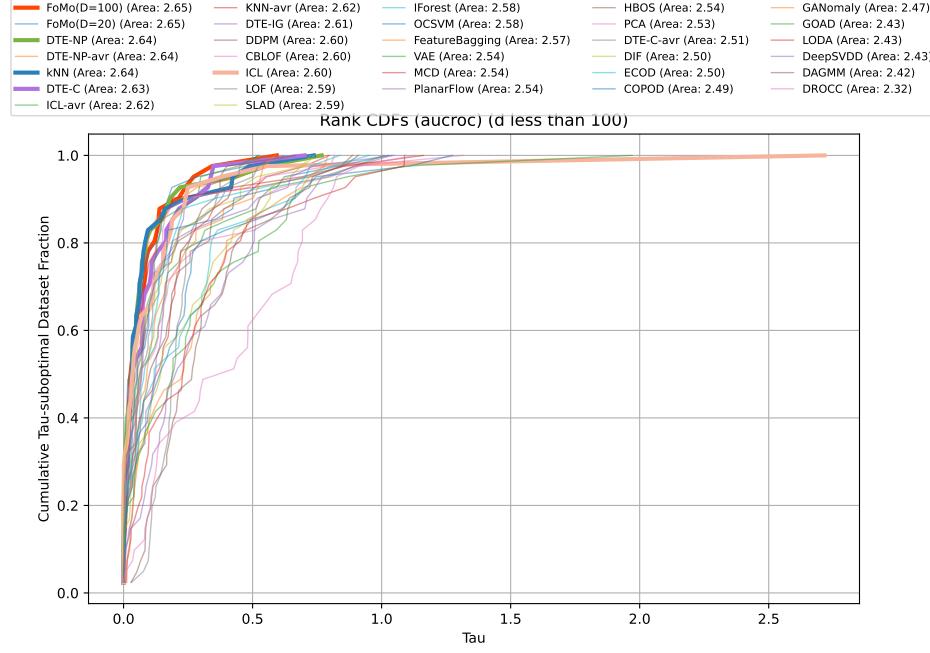


Figure 21: FoMo-0D ($D=100$) ranks at **top-1** based on the performance profile plots of all detectors w.r.t. **AUROC** in *datasets with dimensions less than $d \leq 100$* . In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

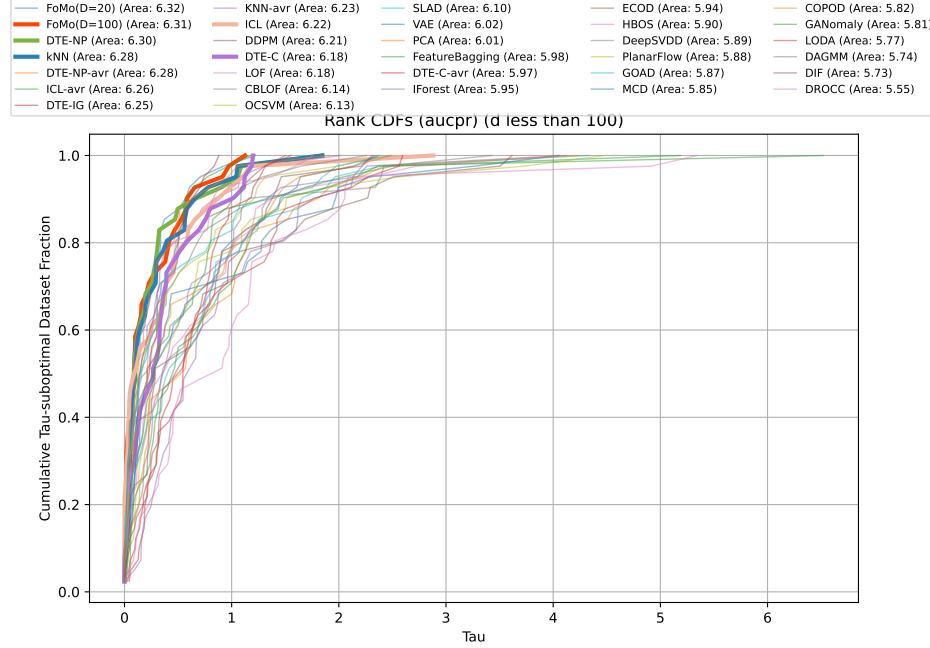


Figure 22: FoMo-0D ($D=100$) ranks at **top-2** based on the performance profile plots of all detectors w.r.t. **AUPR** in *datasets with dimensions less than $d \leq 100$* . In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

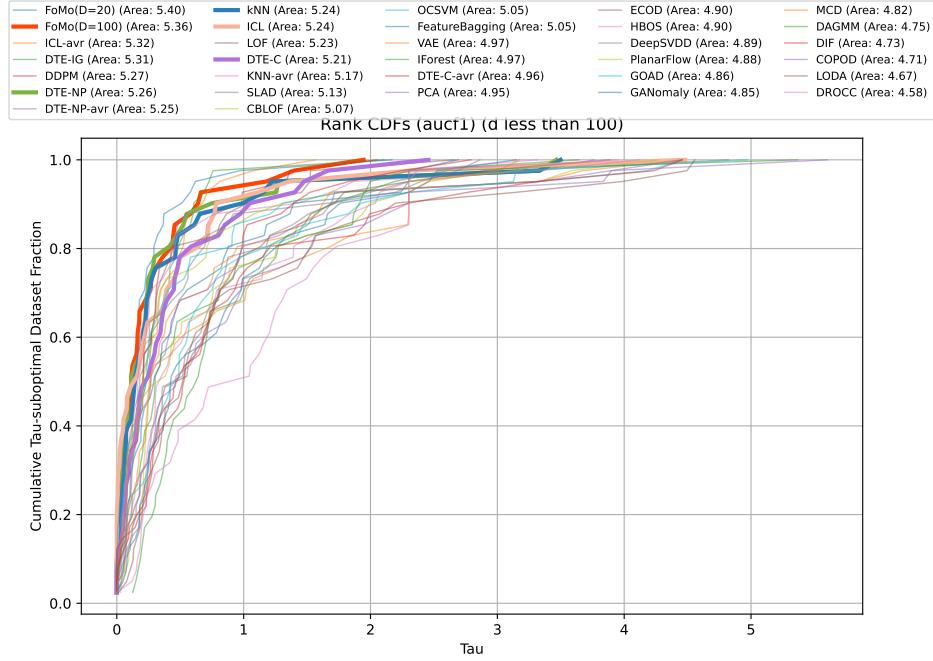


Figure 23: FoMo-0D ($D=100$) ranks at **top-2** based on the performance profile plots of all detectors w.r.t. **F1** in *datasets with dimensions less than $d \leq 100$* . In the plot, x-axis represents the τ values—performance ratios that compare each method’s metric to the best performance achieved, while y-axis displays the cumulative fraction of test datasets for which a method’s performance is within the τ value. We use the area under each CDF curve as a global performance indicator, where a larger area signifies superior performance.

1058 I Full Results

1059 Tables 14.1& 14.2, 15.1 & 15.2, and 16.1 & 16.2 respectively show the AUROC, AUPR, and F1 scores of the
 1060 top-4 baselines, DTE-NP, kNN, ICL, and DTE-C as well as their corresponding ^{avg} model with the average
 1061 performance across HPs, as listed in Table 4.

1062 Tables 17.1&17.2, 18.1&18.2, and 19.1&19.2 respectively show the AUROC, AUPR, and F1 scores of all
 1063 methods across all benchmark datasets. In all these tables, the last four rows show the avg_rank of methods
 1064 across datasets, and p -values of the Wilcoxon signed rank test comparing FoMo-0D w/ $D = 100$ with other
 1065 baselines. The preceding four rows are the same for FoMo-0D w/ $D = 20$, when ranking 31 models (26
 1066 baselines + 4 ^{avg} variants of top-4 baselines + FoMo-0D w/ $D = 20$).

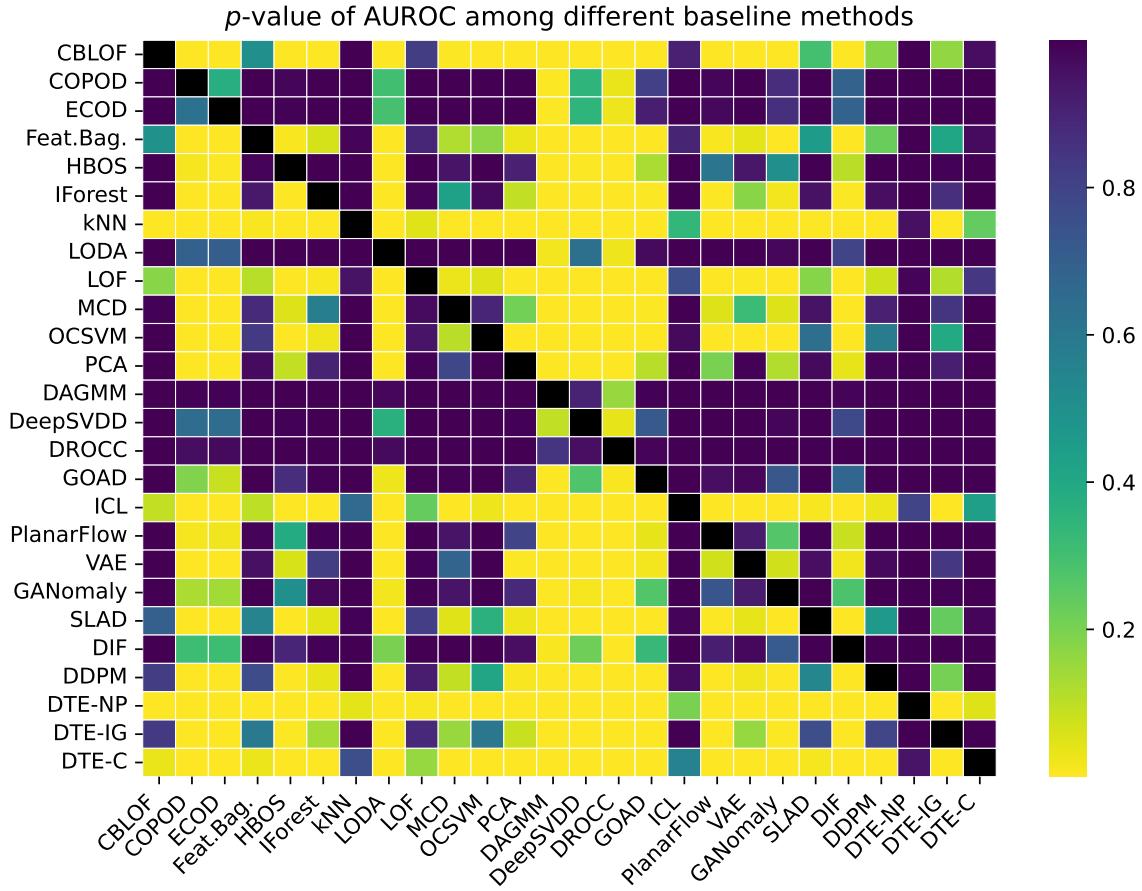


Figure 24: Pairwise *p*-values among baseline methods based on the Wilcoxon signed rank test w.r.t. AUROC performances across datasets.

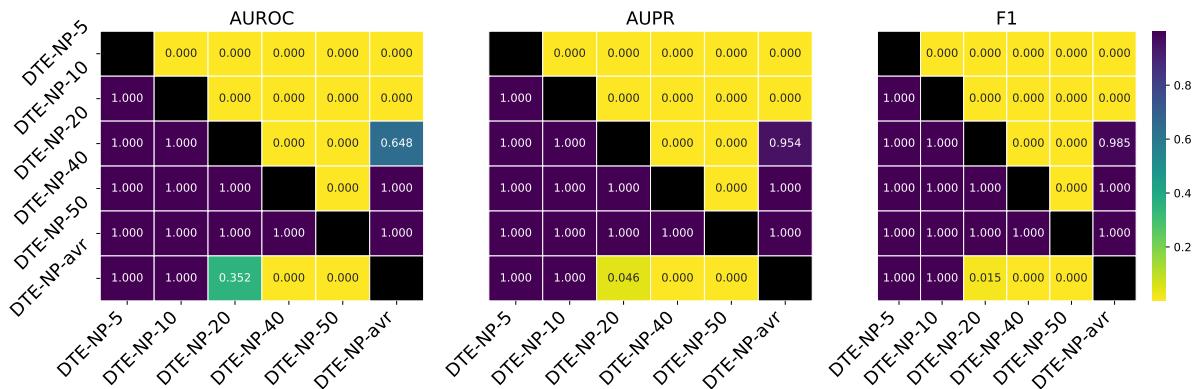


Figure 25: *p*-values w.r.t. AUROC/AUPR/F1 among different HP configurations of **DTE-NP** (i.e., $k \in \{5, 10, 20, 40, 50\}$), along with the avg model with the average performance across HPs.

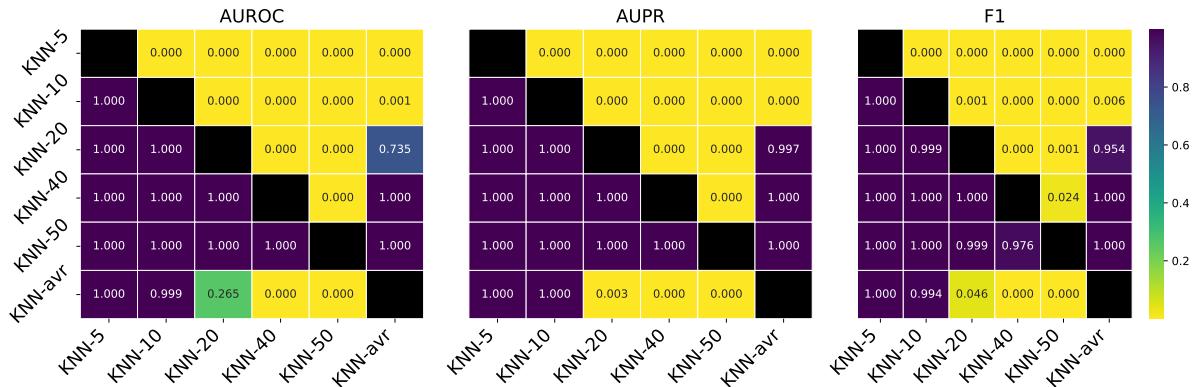


Figure 26: p -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **kNN** (i.e., $k \in \{5, 10, 20, 40, 50\}$), along with the ^{avg} model with the average performance across HPs.

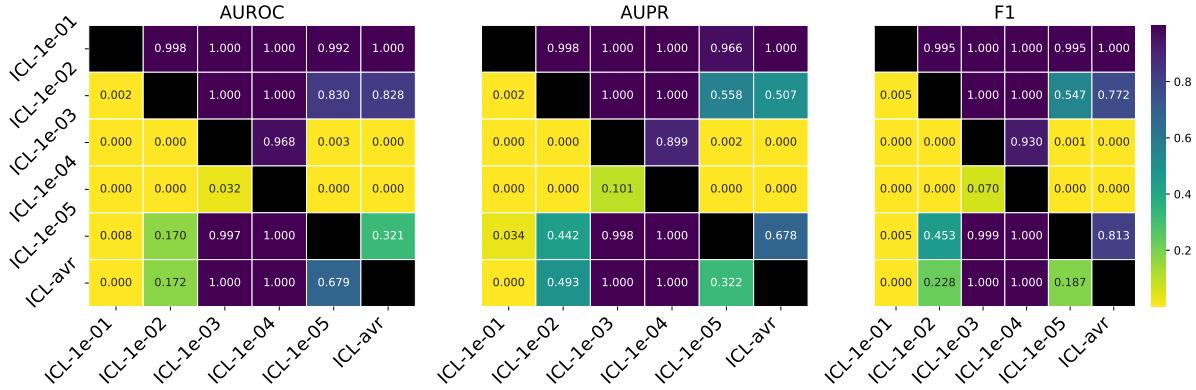


Figure 27: p -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **ICL** (i.e., $\text{learning_rate} \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$), along with the ^{avg} model with the average performance across HPs.

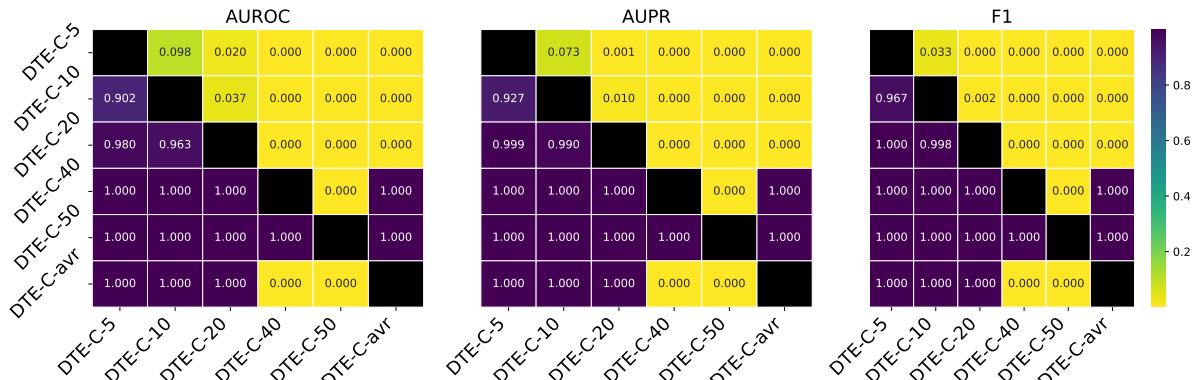


Figure 28: p -values w.r.t. AUROC/AUPR/F1 among different HP configurations of **DTE-C** (i.e., $k \in \{5, 10, 20, 40, 50\}$), along with the ^{avg} model with the average performance across HPs.

Table 13: Comparison of methods across datasets. (top row) Rank w.r.t. AUROC performance avg.'ed over 57 datasets is presented for FoMo-0D (with $D = 100$), **top-10 baselines** with default HPs, and **top-4⁶** baselines with performance **avg.**'ed over varying HPs (denoted w/ ^{avg}); followed by p -values of the pairwise Wilcoxon signed rank test, comparing FoMo-0D to each baseline (from top to bottom) over All (57) datasets, those (24) w/ $d \leq 20$, (38) w/ $d \leq 50$, (42) w/ $d \leq 100$ and (46) datasets w/ $d \leq 500$ dimensions. FoMo-0D performs as well as (**i.e., statistically no different from**) the **2nd best model** (kNN , w/ $p = 0.106$) across All datasets, while it is **comparable to** ($p > 0.05$) or **better than** ($p > 0.95$) **all baselines** over datasets w/ $d \leq 100$ (aligned w/ pretraining where $D = 100$) *and* $d \leq 500$ (generalizing beyond pretraining).

		FoMo-0D	DTE-NP	kNN	ICL	DTE-C	LOF	CBLOF	Feat.Bag.	SLAD	DDPM	OCSVM	DTE-NP ^{avg}	kNN^{avg}	ICL ^{avg}	DTE-C ^{avg}	
Rank(avg)		11.886		7.553	9.018	10.851	11.36	12.316	13.342	13.386	12.982	14.061	13.851	9.079	11.105	12.991	22.263
All	-		<u>0.016</u>	0.106	0.462	0.454	0.585	0.750	0.823	0.759	0.901	0.895		0.112	0.315	0.670	1.000
$d \leq 20$	-		0.428	0.665	0.987	0.727	0.911	0.940	0.987	0.868	0.758	0.968		0.781	0.868	0.990	1.000
$d \leq 50$	-		0.734	0.923	0.992	0.973	0.989	0.987	0.999	0.948	0.985	0.986		0.948	0.967	0.989	1.000
$d \leq 100$	-		0.415	0.700	0.949	0.953	0.970	0.971	0.996	0.876	0.980	0.978		0.752	0.860	0.958	1.000
$d \leq 200$	-		0.315	0.605	0.923	0.919	0.944	0.977	0.990	0.904	0.970	0.983		0.663	0.789	0.937	1.000
$d \leq 500$	-		0.220	0.569	0.827	0.894	0.960	0.968	0.994	0.910	0.960	0.979		0.607	0.756	0.846	1.000

J Benchmark OD Datasets

Table 20: Description of all datasets in ADBench Livernoche et al. (2024). Datasets in blue are image and text datasets that are vectorized through pretrained encoders. We refer to the original paper for details.

Dataset Name	# Samples	# Features	# Anomaly	% Anomaly	Category
ALOI	49534	27	1508	3.04	Image
anthyroid	7200	6	534	7.42	Healthcare
backdoor	95329	196	2329	2.44	Network
breastw	683	9	239	34.99	Healthcare
campaign	41188	62	4640	11.27	Finance
cardio	1831	21	176	9.61	Healthcare
Cardiotocography	2114	21	466	22.04	Healthcare
celeba	202599	39	4547	2.24	Image
census	299285	500	18568	6.20	Sociology
cover	286048	10	2747	0.96	Botany
donors	619326	10	36710	5.93	Sociology
fault	1941	27	673	34.67	Physical
fraud	284807	29	492	0.17	Finance
glass	214	7	9	4.21	Forensic
Hepatitis	80	19	13	16.25	Healthcare
http	567498	3	2211	0.39	Web
InternetAds	1966	1555	368	18.72	Image
Ionosphere	351	32	126	35.90	Oryctognosy
landsat	6435	36	1333	20.71	Astronautics
letter	1600	32	100	6.25	Image
Lymphography	148	18	6	4.05	Healthcare
magic.gamma	19020	10	6688	35.16	Physical
mammography	11183	6	260	2.32	Healthcare
mnist	7603	100	700	9.21	Image
musk	3062	166	97	3.17	Chemistry
optdigits	5216	64	150	2.88	Image
PageBlocks	5393	10	510	9.46	Document
pendigits	6870	16	156	2.27	Image
Pima	768	8	268	34.90	Healthcare
satellite	6435	36	2036	31.64	Astronautics
satimage-2	5803	36	71	1.22	Astronautics
shuttle	49097	9	3511	7.15	Astronautics
skin	245057	3	50859	20.75	Image
smtp	95156	3	30	0.03	Web
SpamBase	4207	57	1679	39.91	Document
speech	3686	400	61	1.65	Linguistics
Stamps	340	9	31	9.12	Document
thyroid	3772	6	93	2.47	Healthcare
vertebral	240	6	30	12.50	Biology
vowels	1456	12	50	3.43	Linguistics
Waveform	3443	21	100	2.90	Physics
WBC	223	9	10	4.48	Healthcare
WDBC	367	30	10	2.72	Healthcare
Wilt	4819	5	257	5.33	Botany
wine	129	13	10	7.75	Chemistry
WPBC	198	33	47	23.74	Healthcare
yeast	1484	8	507	34.16	Biology
CIFAR10	5263	512	263	5.00	Image
FashionMNIST	6315	512	315	5.00	Image
MNIST-C	10000	512	500	5.00	Image
MVTec-AD	5354	512	1258	23.50	Image
SVHN	5208	512	260	5.00	Image
Agnews	10000	768	500	5.00	NLP
Amazon	10000	768	500	5.00	NLP
Imdb	10000	768	500	5.00	NLP
Yelp	10000	768	500	5.00	NLP
20newsgroups	11905	768	591	4.96	NLP

1068 K Differences to Prior Work on PFNs for Tabular Data

1069 There exist applications of PFNs (originally developed by Müller et al. (2022)) that pre-date our proposed
 1070 FoMo-OD, namely, TabPFN (Hollmann et al., 2023) for supervised classification, LC-PFN (Adriaensen
 1071 et al., 2024) for learning curve extrapolation, PFN4BO (Müller et al., 2023) for Bayesian optimization, and
 1072 ForecastPFN (Dooley et al., 2023) for time series forecasting.

1073 Here we highlight the differences of our proposed FoMo-OD from these existing PFNs.

- 1074 1. **First PFN4OD:** We employ prior-data fitted networks (PFNs) for outlier detection (OD) for the
 1075 first time.
- 1076 2. **First large-scale pretrained OD model:** FoMo-OD is the first model for zero-shot OD that is
 1077 pretrained at large scale on a large collection of (synthetic) datasets, due to the minuscule nature of
 1078 existing real-world OD benchmark datasets.
- 1079 3. **New data prior:** Thanks to PFN’s reliance on synthetically generated datasets, we establish a new
 1080 data prior for OD, specifically for outlier synthesis.
- 1081 4. **Data transformation for scale:** While drawing samples from a data prior may be relatively fast,
 1082 pretraining a large foundation model requires many such draws for every step of each epoch. To speed
 1083 up data synthesis on-the-fly, we are the first to leverage a linear transformation.
- 1084 5. **Router-based attention for scale:** PFNs ingest the entire training dataset as context for in-context
 1085 learning at inference time. To accommodate larger datasets at both training (for better generalization)
 1086 and inference (for large-scale real-world datasets), we leveraged a “bottleneck” architecture for scalable
 1087 self-attention, and in turn, larger context size.

1088 L Discussion

1089 **Summary:** We introduced FoMo-OD, **the first foundation model for outlier detection (OD)** on
 1090 tabular data. FoMo-OD is a prior-data fitted network (PFN), pretrained on a large number of *synthetic*
 1091 datasets generated from a new data prior for OD, which can infer the posterior predictive distribution for
 1092 test points in a new dataset in a **zero-shot** fashion where the training data is input as context, capitalizing
 1093 on *in-context learning*.

1094 Zero-shot OD implies **no additional OD model training or model selection**, given a new OD task. That
 1095 is a revolution for OD (!), for which algorithm and hyperparameter selection are notoriously-hard *without any*
 1096 *labeled data*, and also computationally taxing especially for today’s modern deep OD models with numerous
 1097 parameters *and* a long list of hyperparameters. What is more, FoMo-OD provides **extremely fast inference**
 1098 thanks to a mere *single forward pass*, making it amenable for OD on data streams.

1099 Building on the PFN paradigm (Müller et al., 2022), FoMo-OD breaks new ground not only conceptually
 1100 by abolishing the burden of model training and selection, but also empirically: Against **26** different (both
 1101 classical and modern) baselines on **57** public benchmark datasets from diverse domains, FoMo-OD performs on
 1102 par with the top **2nd** baseline, while significantly outperforming the majority of the baselines. Without the
 1103 need to train any, let alone multiple models for HP tuning, FoMo-OD takes a mere **7.7 ms** per test sample for
 1104 inference only.

1105 **Limitations and Future Directions:** FoMo-OD employs a simple straightforward data prior based on
 1106 GMMs. While it is remarkable to see how far one can go with synthetic data from such a simple prior, future
 1107 work can design more comprehensive data priors, inclusive of discrete features as well as other possible outlier
 1108 types. We have also pretrained FoMo-OD solely on synthetic datasets, while future work can augment both
 1109 synthetic and real-world datasets for pretraining.

1110 Besides the lack of massive real-world datasets for tabular OD, a motivation for a data prior to pretrain purely
 1111 on synthetic datasets comes from neural scaling laws (Kaplan et al., 2020; Zhai et al., 2022). Interestingly,
 1112 the scaling laws for large Transformer models have shown that their generalization error tends to drop as a
 1113 power law with the amount of training data (also, with number of parameters and amount of compute), but

1114 the power law exponent is very small—suggesting that acquiring more colossal real-world datasets would be a
 1115 slow, if not expensive approach to advancing ML/AI. Others have proposed ways to subset-select smaller,
 1116 non-redundant “foundation datasets” (Sorscher et al., 2022; Paul et al., 2021), and emphasized the importance
 1117 of task/dataset diversity in pretraining (Raventós et al., 2024). Arguably, synthetic data from a complex and
 1118 diverse data prior is a potential gateway to obtaining non-redundant and diverse datasets for pretraining
 1119 large foundation models like FoMo-0D. On the other hand, designing such a data prior requires a level of
 1120 domain/prior knowledge.

1121 Another improvement could be scaling up to even larger context (i.e. dataset) size and dimensionality.
 1122 While FoMo-0D generalizes beyond pretrained context sizes and dimensionality, it is limited to and performs
 1123 particularly well on downstream datasets of similar nature as our experiments showed. A promising direction
 1124 for size generalization is using PFNs as extremely fast ensemble components at inference; since “*PFNs are*
 1125 *quick enough to be used as ensemble members. The size constraints could therefore be overcome by boosting*
 1126 *and bagging techniques*” (Nagler, 2023).

1127 Further, our work focused on unsupervised OD with clean/inlier-only training data. Future work can study
 1128 the unsupervised OD setting and pretraining with mixed/“contaminated” data in the transductive setting,
 1129 where the unlabeled test data is the same as training data. In addition, we performed offline evaluation
 1130 of FoMo-0D on static datasets, while its fast inference lends itself to streaming OD, which future work can
 1131 explore. Technically, both extensions (unsupervised OD and streaming OD) are straightforward from the
 1132 implementation perspective.

1133 Our current work is limited to OD for tabular (or point-cloud) data. Our ideas can be extended to other data
 1134 modalities, such as image, graph, and text outliers, to comprise other domains with critical OD applications
 1135 such as video surveillance, fraud detection and LLM hallucination detection. To that end, the design of
 1136 novel inlier/outlier priors would be an open direction. A promising approach here could be the use of
 1137 pretrained generative models to draw synthesized image/text/etc. datasets for pretraining the PFN, in place
 1138 of manually-designed data priors.

1139 Finally, our quest here has been mainly experimental. Theoretically understanding why these models work as
 1140 well as they do and investigating their failure cases are important yet open questions. *As empirical future*
 1141 *work, one could systematically stress-test FoMo-0D using synthetically generated test datasets that contain*
 1142 *known outlier types and inlier distributions distinct from those used during pretraining, while varying the*
 1143 *degree of out-of-distribution characteristics in a controlled manner.*

1144 As the first foundation model for OD, FoMo-0D inspires many promising directions for future research that
 1145 could lead to fruition for additional practical applications.

1146 M Reproducibility Statement

1147 We expect that the disruptive nature of FoMo-0D will trigger future innovations in the OD literature, as
 1148 well as a widespread adoption by practitioners thanks to its key desirable properties. To foster future
 1149 research and accessibility in practice, we make all resources (our codebase used for prior data synthesis, data
 1150 transformation, and pretraining as well as our pretrained model checkpoints) publicly available at <https://anonymous.4open.science/r/PFN40D>. Further, full implementation details are provided in Appendix C.